

PCT

WORLD INTELLECTUAL PROPERTY ORGANIZATION  
International Bureau



B32

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

<p>(51) International Patent Classification :  Not classified</p>	<p>A2</p>	<p>(11) International Publication Number: <b>WO 98/58529</b>  (43) International Publication Date: 30 December 1998 (30.12.98)</p>
<p>(21) International Application Number: PCT/US98/12930 (22) International Filing Date: 22 June 1998 (22.06.98)  (30) Priority Data: 60/050,594 24 June 1997 (24.06.97) US  (63) Related by Continuation (CON) or Continuation-in-Part (CIP) to Earlier Application US 60/050,594 (CON) Filed on 24 June 1997 (24.06.97)  (71) Applicant (for all designated States except US): AFFYMETRIX, INC. [US/US]; 3380 Central Expressway, Santa Clara, CA 95051 (US).  (72) Inventors; and (75) Inventors/Applicants (for US only): LIPSHUTZ, Robert, J. [US/US]; 970 Palo Alto Avenue, Palo Alto, CA 94301 (US). CHEE, Mark [AU/US]; 3199 Waverly Street, Palo Alto, CA 94306 (US). FAN, Jian-Bing [CN/US]; Apartment 20, 275 Ventura Avenue, Palo Alto, CA 94306 (US). BERNO, Anthony [CA/US]; 570 South 12th Street, San Jose, CA 95112 (US).</p>		<p>(74) Agents: LIEBESCHUETZ, Joe et al.; Townsend and Townsend and Crew LLP, 8th floor, Two Embarcadero Center, San Francisco, CA 94111-3834 (US).  (81) Designated States: JP, US, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).  <b>Published</b> <i>With declaration under Article 17(2)(a); without classification and without abstract; title not checked by the International Searching Authority.</i></p>
<p>(54) Title: GENETIC COMPOSITIONS AND METHODS</p>		

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

## GENETIC COMPOSITIONS AND METHODS

## BACKGROUND OF THE INVENTION

The genomes of all organisms undergo spontaneous mutation in the course of their continuing evolution generating variant forms of progenitor sequences (Gusella, *Ann. Rev. Biochem.* 55, 831-854 (1986)). The variant form may confer an evolutionary advantage or disadvantage relative to a progenitor form or may be neutral. In some instances, a variant form confers a lethal disadvantage and is not transmitted to subsequent generations of the organism. In other instances, a variant form confers an evolutionary advantage to the species and is eventually incorporated into the DNA of many or most members of the species and effectively becomes the progenitor form. In many instances, both progenitor and variant form(s) survive and co-exist in a species population. The coexistence of multiple forms of a sequence gives rise to polymorphisms.

Several different types of polymorphism have been reported. A restriction fragment length polymorphism (RFLP) means a variation in DNA sequence that alters the length of a restriction fragment as described in Botstein et al., *Am. J. Hum. Genet.* 32, 314-331 (1980). The restriction fragment length polymorphism may create or delete a restriction site, thus changing the length of the restriction fragment. RFLPs have been widely used in human and animal genetic analyses (see WO 90/13668; W090/11369; Donis-Keller, *Cell* 51, 319-337 (1987); Lander et al., *Genetics* 121, 85-99 (1989)). When a heritable trait can be linked to a particular RFLP, the presence of the RFLP in an individual can be used to predict the likelihood that the animal will also exhibit the trait.

Other polymorphisms take the form of short tandem repeats (STRs) that include tandem di-, tri- and tetra-nucleotide repeated motifs. These tandem repeats are also referred to as variable number tandem repeat (VNTR) polymorphisms. VNTRs have been used in identity and paternity analysis (US 5,075,217; Armour et al., *FEBS Lett.* 307, 113-115 (1992); Horn et al., WO 91/14003; Jeffreys, EP 370,719), and in a large number of genetic mapping studies.

Other polymorphisms take the form of single nucleotide variations between individuals of the same species. Such polymorphisms are far more frequent than RFLPs, STRs and VNTRs. Some single nucleotide polymorphisms occur in protein-coding sequences, in which case, one of the polymorphic forms may give rise to the expression of a defective or other variant protein and, potentially, a genetic disease. Examples of genes, in which polymorphisms within coding sequences give rise to genetic disease include  $\beta$ -globin (sickle cell anemia) and CFTR (cystic fibrosis). Other single nucleotide polymorphisms occur in noncoding regions. Some of these polymorphisms may also result in defective protein expression (e.g., as a result of defective splicing). Other single nucleotide polymorphisms have no phenotypic effects.

Single nucleotide polymorphisms can be used in the same manner as RFLPs, and VNTRs but offer several advantages. Single nucleotide polymorphisms occur with greater frequency and are spaced more uniformly throughout the genome than other forms of polymorphism. The greater frequency and uniformity of single nucleotide polymorphisms means that there is a greater probability that such a polymorphism will be found in close proximity to a genetic locus of interest than would be the case for other polymorphisms. Also, the different forms of characterized single nucleotide polymorphisms are often easier to distinguish than other types of polymorphism (e.g., by use of assays employing allele-specific hybridization probes or primers).

Despite the increased amount of nucleotide sequence data being generated in recent years, only a minute proportion of the total repository of polymorphisms in humans and other

organisms has so far been identified. The paucity of polymorphisms hitherto identified is due to the large amount of work required for their detection by conventional methods. For example, a conventional approach to identifying polymorphisms might be to sequence the same stretch of oligonucleotides in a population of individuals by dideoxy sequencing. In this type of approach, the amount of work increases in proportion to both the length of sequence and the number of individuals in a population and becomes impractical for large stretches of DNA or large numbers of persons.

#### SUMMARY OF THE CLAIMED INVENTION

The invention provides nucleic acid segments of between 10 and 100 bases from a fragment shown in Table 1, column 1 including a polymorphic site. Complements of these segments are also included. The segments can be DNA or RNA, and can be double- or single-stranded. Some segments are 10-20 or 10-50 bases long. Preferred segments include a diallelic polymorphic site. The base occupying the polymorphic site in the segments can be the reference (Table 1, column 3) or an alternative base (Table 1, column 5).

The invention further provides allele-specific oligonucleotides that hybridizes to a segment of a fragment shown in Table 1, column 8 or its complement. These oligonucleotides can be probes or primers. Also provided are isolated nucleic acids comprising a sequence of Table 1, column 8, or the complement thereto, in which the polymorphic site within the sequence is occupied by a base other than the reference base shown in Table 1, column 3.

The invention further provides a method of analyzing a nucleic acid from an individual. The method determines which base is present at any one of the polymorphic sites shown in Table 1. Optionally, a set of bases occupying a set of the polymorphic sites shown in Table 1 is determined. This type of analysis can be performed on a plurality of individuals who are tested for the presence of a disease phenotype. The presence or absence of disease phenotype can then be

correlated with a base or set of bases present at the polymorphic sites in the individuals tested.

5 The invention further provides computer-readable storage medium for storing data for access by an application program being executed on a data processing system. Such a medium comprises a data structure stored in the computer-readable storage medium, the data structure including information resident in a database used by the application program. The data structure includes a plurality of records, each record of the plurality comprising information  
10 identifying a polymorphisms shown in Table 1.

The invention further provides a signal carrying data for access by an application program being executed on a data processing system. A data structure is encoded in the signal.  
15 The data structure includes information resident in a database used by the application program. Such information includes a plurality of records, each record of the plurality comprising information identifying a polymorphism shown in Table 1.

#### BRIEF DESCRIPTION OF THE FIGURES

20 Figs. 1A and 1B depict computer systems suitable for storing and transmitting information relating to the polymorphisms of the invention.

#### DEFINITIONS

An oligonucleotide can be DNA or RNA, and single- or  
25 double-stranded. Oligonucleotides can be naturally occurring or synthetic, but are typically prepared by synthetic means. Preferred oligonucleotides of the invention include segments of DNA, or their complements including any one of the polymorphic sites shown in Table 1. The segments are usually  
30 between 5 and 100 bases, and often between 5-10, 5-20, 10-20, 10-50, 15-50, 15-100, 20-50 or 20-100 bases. The polymorphic site can occur within any position of the segment. The segments can be from any of the allelic forms of DNA shown in Table 1.

35 Hybridization probes are oligonucleotides capable of binding in a base-specific manner to a complementary strand of

nucleic acid. Such probes include peptide nucleic acids, as described in Nielsen et al., *Science* 254, 1497-1500 (1991).

The term primer refers to a single-stranded oligonucleotide capable of acting as a point of initiation of template-directed DNA synthesis under appropriate conditions (i.e., in the presence of four different nucleoside triphosphates and an agent for polymerization, such as, DNA or RNA polymerase or reverse transcriptase) in an appropriate buffer and at a suitable temperature. The appropriate length of a primer depends on the intended use of the primer but typically ranges from 15 to 30 nucleotides. Short primer molecules generally require cooler temperatures to form sufficiently stable hybrid complexes with the template. A primer need not reflect the exact sequence of the template but must be sufficiently complementary to hybridize with a template. The term primer site refers to the area of the target DNA to which a primer hybridizes. The term primer pair means a set of primers including a 5' upstream primer that hybridizes with the 5' end of the DNA sequence to be amplified and a 3', downstream primer that hybridizes with the complement of the 3' end of the sequence to be amplified.

Linkage describes the tendency of genes, alleles, loci or genetic markers to be inherited together as a result of their location on the same chromosome, and can be measured by percent recombination between the two genes, alleles, loci or genetic markers.

Polymorphism refers to the occurrence of two or more genetically determined alternative sequences or alleles in a population. A polymorphic marker or site is the locus at which divergence occurs. Preferred markers have at least two alleles, each occurring at frequency of greater than 1%, and more preferably greater than 10% or 20% of a selected population. A polymorphic locus may be as small as one base pair. Polymorphic markers include restriction fragment length polymorphisms, variable number of tandem repeats (VNTR's), hypervariable regions, minisatellites, dinucleotide repeats, trinucleotide repeats, tetranucleotide repeats, simple sequence repeats, and insertion elements such as Alu. The

first identified allelic form is arbitrarily designated as a the reference form and other allelic forms are designated as alternative or variant alleles. The allelic form occurring most frequently in a selected population is sometimes referred to as the wildtype form. Diploid organisms may be homozygous or heterozygous for allelic forms. A diallelic polymorphism has two forms. A triallelic polymorphism has three forms.

A single nucleotide polymorphism occurs at a polymorphic site occupied by a single nucleotide, which is the site of variation between allelic sequences. The site is usually preceded by and followed by highly conserved sequences of the allele (e.g., sequences that vary in less than 1/100 or 1/1000 members of the populations).

A single nucleotide polymorphism usually arises due to substitution of one nucleotide for another at the polymorphic site. A transition is the replacement of one purine by another purine or one pyrimidine by another pyrimidine. A transversion is the replacement of a purine by a pyrimidine or vice versa. Single nucleotide polymorphisms can also arise from a deletion of a nucleotide or an insertion of a nucleotide relative to a reference allele.

Hybridizations are usually performed under stringent conditions, for example, at a salt concentration of no more than 1 M and a temperature of at least 25°C. For example, conditions of 5X SSPE (750 mM NaCl, 50 mM NaPhosphate, 5 mM EDTA, pH 7.4) and a temperature of 25-30°C are suitable for allele-specific probe hybridizations.

An isolated nucleic acid means an object species invention that is the predominant species present (i.e., on a molar basis it is more abundant than any other individual species in the composition). Preferably, an isolated nucleic acid comprises at least about 50, 80 or 90 percent (on a molar basis) of all macromolecular species present. Most preferably, the object species is purified to essential homogeneity (contaminant species cannot be detected in the composition by conventional detection methods).

Linkage disequilibrium or allelic association means the preferential association of a particular allele or genetic



marker with a specific allele, or genetic marker at a nearby chromosomal location more frequently than expected by chance for any particular allele frequency in the population. For example, if locus X has alleles a and b, which occur equally frequently, and linked locus Y has alleles c and d, which occur equally frequently, one would expect the combination ac to occur with a frequency of 0.25. If ac occurs more frequently, then alleles a and c are in linkage disequilibrium. Linkage disequilibrium may result from natural selection of certain combination of alleles or because an allele has been introduced into a population too recently to have reached equilibrium with linked alleles.

A marker in linkage disequilibrium can be particularly useful in detecting susceptibility to disease (or other phenotype) notwithstanding that the marker does not cause the disease. For example, a marker (X) that is not itself a causative element of a disease, but which is in linkage disequilibrium with a gene (including regulatory sequences) (Y) that is a causative element of a phenotype, can be used detected to indicate susceptibility to the disease in circumstances in which the gene Y may not have been identified or may not be readily detectable.

The present invention includes the use of any of the polymorphic forms shown in Table 1 as a means to determine susceptibility to a phenotype resulting from an allele or marker in linkage disequilibrium with such polymorphic forms.

#### DESCRIPTION

##### I. Novel Polymorphisms of the Invention

The novel polymorphisms of the invention are listed in Table 1. The first column of the Table lists the names assigned to the fragments in which the polymorphisms occur. The fragments are all human genomic fragments. SGC, TIGR and WI respectively stand for Stanford Genome Center, The Institute for Genome Research and the Whitehead Institute. The sequence of one allelic form of each of the fragments (arbitrarily referred to as the prototypical or reference

form) has been previously published. These sequences are listed at <http://www-genome.wi.mit.edu/> (all STS's (sequence tag sites)); <http://shgc.stanford.edu> (Stanford STS's); and <http://ww.tigr.org/> (TIGR STS's). The Web sites also list  
5 primers for amplification of the fragments, and the genomic location of fragments. Some fragments are expressed sequence tags, and some are random genomic fragments. All information in the websites concerning the fragments listed in Table 1 is incorporated by reference in its entirety for all purposes.

10 The second column lists the position in the fragment in which a polymorphic site has been found. Positions are numbered consecutively with the first base of the fragment sequence as listed in one of the above databases being assigned the number one. The third column lists the base  
15 occupying the polymorphic site in the sequence in the data base. This base is arbitrarily designated the reference or prototypical form but is not necessarily the most frequently occurring form. The fifth column in the table lists the alternative base(s) at the polymorphic site. The eighth  
20 column of the Table lists about 15 bases of sequence on either side of the polymorphic site in each fragment. The indicated sequences can be either DNA or RNA. In the latter, the T's shown in the Table are replaced by U's. The base occupying the polymorphic site is indicated in EUPAC-IUB ambiguity code.  
25 The fourth and sixth columns of the table show the frequency with which reference and alternative alleles occur at a polymorphic site. The seventh column in the table indicates the population frequency of heterozygotes of the polymorphic site. Also provided is a nucleic acid encoding hepatic lipase  
30 containing a polymorphism. The sequence is  
CTTCGAGAGAGATTGMACAGATTCCTGGAAG.

Table 1

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-16260	59	G	0.79	T	0.21	0.34	GATTCAAGAAAGAAAACCCAGAGTTTCACA
WI-16260	86	G	0.79	A	0.21	0.34	CACAATATAGGTAGCRATAACCAAGGTCTCAC
WI-16303	65	A	0.93	G	0.07	0.13	GGTCACTGCAGCCCCCTCTGTATTAGGGAGC
WI-16398	90	T	0.36	C	0.64	0.46	TCCATGATATTTTCAYAGCAACTAGTATATA
WI-16403	69	T	0.71	C	0.29	0.41	TCAATTTTAAACACTYCTTTTATATAGGGA
WI-16406	24	C	0.86	T	0.14	0.24	GCTACAGAAAGAGGYGGTTTATTTTCTTT
WI-16543	67	G	0.50	T	0.50	0.50	ACATTTGGGTTTIGGKAAAGTCCCTGTAAATG
WI-16632	71	A	0.44	G	0.56	0.49	CTACTTTGGAGCCCTRAGGAGTTTITAGAGA
WI-16644	42	G	0.25	A	0.75	0.38	GCTCATTTTGATTACRGGGTATACATGAAGTA
WI-16739	57	G	0.44	A	0.56	0.49	TTTGCCATCACAAAGCCTTATAGGGAATAATG
WI-16782	96	C	0.69	T	0.31	0.43	GTCTCACTGTAAAGGAYGATGGAGGAACAGAA
WI-16783	64	A	0.75	G	0.25	0.38	TGTCCTTTACCTGAGRCTAATAAGGATTGAA
WI-16816	124	A	0.75	G	0.25	0.38	CCATTGTTGGGGTTARACTGTCTGAACAAA
WI-16824	47	T	0.25	C	0.75	0.38	TGGTGTGCAGCTGTGTCTTCTTATGAAGAAG
WI-16824	83	G	0.75	A	0.25	0.38	AGCTGATAAACGTGGRCTTACACCTTTAGCA
WI-16857	47	G	0.13	A	0.88	0.22	GCAGCTAAATGGCAAATRCTAGTGGTCTTCCCA
WI-16879	79	C	0.88	T	0.13	0.22	AGGCCATATTTTCCCAATATAGGACTCTAGTTC

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-16882	99	A	0.56	G	0.44	0.49	TGCCACGTCTCTGACRGCAGATTACCTGACA
WI-16888	70	G	0.38	A	0.63	0.47	ACTTTGGGCAGGTTCTRTTAAATTTGGTCAAT
WI-16905	75	C	0.88	T	0.13	0.22	GGCCTGTGTTGTTTCAYCCCACTGCCCTAGAAG
WI-16910	74	G	0.75	A	0.25	0.38	AAGATGGCGCTAGAAAGTATCTGTTATAGAA
WI-16918	93	C	0.44	T	0.56	0.49	CATTAACACCAGCACYCATGCCACTTCTGTA
WI-16947	58	C	0.31	G	0.69	0.43	GAAATAGGCCTGGAGSACAGGATTTGGCTGA
WI-16947	127	A	0.38	C	0.63	0.47	AAAGCAGACCTGGGMCACGGGCAATCACA
WI-16966	43	T	0.88	C	0.13	0.22	CATAACAACCTAATAAYCTTAACTTGGTCCAA
WI-16992	46	G	0.38	A	0.63	0.47	CAGAAGTACACTGTCTCCCTCATCTGAGATG
WI-16992	60	T	0.44	G	0.56	0.49	CGCCCTCATCTGAGAKGTGTAGGACTGTAAAG
WI-16995	55	T	0.25	C	0.75	0.38	GAGGTAAATAGTATTYACGGCTGGAAATCAA
WI-17010	23	T	0.81	C	0.19	0.30	ACAGGAAAAGCCATGYATGACATTCAAAACA
WI-17021	62	T	0.88	A	0.13	0.22	AGCCTATAACTACTCTCWCAGCTGCCACTAAC
WI-17040	94	T	0.44	C	0.56	0.49	ATCATCTCAAGCCAGYCATCACTGAATAAGC
WI-17044	47	G	0.69	T	0.31	0.43	GGATTAACGTATAGGKTCTTAAACAAGGGGA
WI-17065	90	T	0.31	C	0.69	0.43	GAAAAGCATAAACTTYAGGATTTTCATTGCT
WI-17066	32	A	0.38	C	0.63	0.47	CCAACATCACTGTTTMAATTCAGAACATTTT
WI-17074	86	T	0.94	G	0.06	0.12	CTCCTACACAGGCCTKCTACATAGGAGTATA
WI-17104	108	T	0.88	C	0.13	0.22	GGTTTCCAGACGGGCTYTCTCTTTGTAAAGAA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-17108	74	C	0.81	T	0.19	0.30	TCTCAAAGTAAACACYGGGAGCATATGATAA
WI-17114	37	T	0.44	C	0.56	0.49	CAAGGACTTTGTTTTYGTCTCTTCACTCTGC
WI-17136	33	C	0.94	G	0.06	0.12	ATGTCCTAAATGTSATTC AACATATATGC
WI-17149	48	C	0.44	G	0.56	0.49	TTGAAGGAGGAACATSTCATGCACGTGCGTG
WI-17149	79	T	0.31	C	0.69	0.43	GAAACCCAATTGTCA YGTGTATGAACTACAA
WI-17150	76	T	0.38	G	0.63	0.47	GATAGTCTTCTCTCTTKCATATCTTCCAGGAT
WI-17156	54	G	0.81	C	0.19	0.30	TTAGATATCTCCCATSTTCCACAGAAATCAAA
WI-17163	43	A	0.75	G	0.25	0.38	AATAACAATAACGTTTAAAGGCAAAAAGCAAGA
WI-17177	23	A	0.94	G	0.06	0.12	CATATCCAAACCAACCRITCCATCCCCACCTGT
WI-17178	127	T	0.88	C	0.13	0.22	TCCCTCATGAGGAGCYAGAAAGCAGTTGAAAA
WI-17180	47	T	0.75	C	0.25	0.38	AGAGAAATCCTGCACTYCCCAAGTCTCGTCGC
WI-17180	81	C	0.94	G	0.06	0.12	GGCTTCAACAATTACSAACATCTTGCCCCATT
WI-17197	67	G	0.56	A	0.44	0.49	AGTAGCTGGGGCTACRGGTATGCACCACTC
WI-17198	38	A	0.75	C	0.25	0.38	CCTTGTCCTAGTTTMTAAITTTCTCAGTGGA
WI-17347	50	A	0.25	G	0.75	0.38	AGAACTTCTCAGCCTRTGTAGCACAAAGTGGAT
WI-17387	55	C	0.81	G	0.19	0.30	CAGATTGAAGAAAAAATAATAGTAGTTAC
WI-17470	83	A	0.69	G	0.31	0.43	CGTCCCGCCAGCCCTRTCGGCCCTCGTCACTG
WI-17519	55	T	0.38	C	0.63	0.47	TAGCTAATGAATGCAYAGAGTATTGCCTGCA
WI-17581	86	T	0.13	C	0.88	0.22	CCAGTTATTGATAAAYGATAGAACCCAACTA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-17581	99	C	0.69	T	0.31	0.43	AATGATAGAACCCAAATAGGCGCAATTACAA
WI-17596	86	A	0.63	G	0.38	0.47	TGTGTAAACACTCCCRATATTGTGCAATTCT
WI-17623	46	T	0.94	C	0.06	0.12	AATGGTGGGCACATTYGCATGTGCTTACTGG
WI-17675	103	T	0.44	C	0.56	0.49	TTTGATGGTGACTTYCCTGGGTGGTTCCCC
WI-17687	107	C	0.81	G	0.19	0.30	AAAAAGGTGGGGAASGTGCTGGTCGGTACAA
WI-17690	63	G	0.69	A	0.31	0.43	TTTCTAGCTGTGTTTTRATTTGGCTTCCCTAT
WI-17690	79	A	0.63	G	0.38	0.47	ATTGGCTTCCCTATRGATTTCAGGACCCATA
WI-17724	50	T	0.81	C	0.19	0.30	TGGGCCCTCCCTGTCTCYGGACACTGCCAAACCC
WI-17730	39	A	0.44	C	0.56	0.49	AAGTGAAAGTGCTATTMGTTACATCATACCAA
WI-17730	68	T	0.94	C	0.06	0.12	AAGGTACATACTGTGCACATGATTTATGGC
WI-17800	29	C	0.88	G	0.13	0.22	CAAGAGAAACTCACTSAAGACTGGGATTAAT
WI-17835	30	G	0.38	A	0.63	0.47	TATTGTGCTTCTTGRGCCTGTTTCCTATAC
WI-17857	34	T	0.44	G	0.56	0.49	CTGGGATGACTTTCCKATTCTACATCAAGTA
WI-17860	121	T	0.81	A	0.19	0.30	CCAGCAAAGCAAATAWCCGACTGACTGCTCC
WI-17866	43	A	0.63	T	0.38	0.47	CTTCTCAAAATTGTTTWTGTTGTGATTAGTG
WI-17892	76	T	0.88	C	0.13	0.22	GTTTGAGATCACATAYCTGTCTCAGTAGTCT
WI-17904	50	A	0.31	G	0.69	0.43	CAATAAAATGAACACRTACGGGAATTACTAT
WI-17982	98	C	0.25	T	0.75	0.38	ATAACTCCTAAAAGCYGGAAGGAGTTATTAT
WI-17993	118	A	0.94	C	0.06	0.12	CTGTCCCTGTAATGTMCTGCTGAGAGTCCAC

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-17996	84	A	0.13	G	0.88	0.22	AGGCGAAGGGAACAAGRGCTGCCCATGTGCCCT
WI-18012	22	T	0.44	A	0.56	0.49	TGGGTCAGCTCCTTCWTAATGGCCTGAAGGT
WI-18012	46	T	0.38	C	0.63	0.47	TGAAGGTCATCTCCTYTCAACTTCCAGACT
WI-18012	112	C	0.50	T	0.50	0.50	CCACTTTTGCCCCCTTYGTGAAGTGTTCCTG
WI-18012	113	G	0.56	A	0.44	0.49	CACTTTGGCCCCCTTCRTGAAAGTGTTCCTGA
WI-18012	117	A	0.31	G	0.69	0.43	TTTGCCCCCTTCGTGARGTGTTCCTGATACA
WI-18041	24	A	0.75	C	0.25	0.38	AAAAGGTGCTCTCCMGTTTCTAACTCCCTG
WI-18052	50	T	0.31	C	0.69	0.43	TTTCATGTACGAATCYTGGTTACACATCTTA
WI-18052	67	A	0.31	G	0.69	0.43	GGTTACACATCTTAGRACAGCAGAGCTGCCT
WI-18054	46	G	0.25	A	0.75	0.38	AGTGGGGAGTAAARTGGAGCAGGGTGAC
WI-18064	54	G	0.81	A	0.19	0.30	AAGCTGTATTTTCAGARGAATGTCACAATCAT
WI-18068	89	G	0.94	C	0.06	0.12	ATAAAGTAAGACCASATAAAAATACCTATG
WI-18070	28	A	0.88	C	0.13	0.22	ACTCAGAGTGTGTATMATATTAACACATGAA
WI-18080	41	T	0.19	C	0.81	0.30	TCAAACCTAGTCTCTCYTTGTAATTAATAATCT
WI-18080	65	G	0.38	A	0.63	0.47	AAATCTACTATGCCRTGTTTGACTTTTATC
WI-18086	63	G	0.06	A	0.94	0.12	AGAAAGCATACTTCRTGGCTTTGTTACACG
WI-18115	70	C	0.88	T	0.13	0.22	CTTTGGTATCCCTTYCTTTGGTATGAAAGA
WI-18115	71	C	0.88	T	0.13	0.22	TTTGGTATCCCTTCYTTTGGTATGAAAGAC
WI-18136	78	A	0.94	G	0.06	0.12	CTTTAGGTAATTTGCRATAAGAACATAAAAG

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-18169	115	A	0.44	G	0.56	0.49	TCCTCCGGGAAGCTCTGTGGAGCACAAAGCAGA
WI-18181	100	A	0.63	C	0.38	0.47	CACCTCCCTTCAGATCMCAAAAAGCTTAACAAA
WI-18190	62	G	0.75	A	0.25	0.38	GAAGCTAATCATGGARGCAAGCTCCCTGGAG
WI-18215	78	G	0.75	A	0.25	0.38	AGAGTTCCTGCCCTCRGTGTGCGGGGGGAGA
WI-18232	60	T	0.75	A	0.25	0.38	TGTGATACACTTAAGWGAACCCCTGAAAAACC
WI-18242	30	G	0.88	A	0.13	0.22	TAATCGTAACATACTRGAAAGCTGTTACAGT
WI-18266	97	C	0.38	T	0.63	0.47	TGGACTATCTTCAAAATGACAAATGATGCA
WI-18266	124	T	0.13	C	0.88	0.22	TGCATGAATCCACATYTGAGACCCCGCAACTC
WI-18312	73	A	0.75	G	0.25	0.38	ATTGTTATTTCAAAATRTATCTTCTGCTCCCT
WI-18327	104	G	0.44	A	0.56	0.49	TTCGTTAGGCTAGTTRGCTGAGCCATTGTAT
WI-18330	49	G	0.63	A	0.38	0.47	AAATCAGGGATAAGARCTGAGGAACAAGAGG
WI-18357	89	C	0.63	G	0.38	0.47	AGCCCTTAGCATCAASTCATCTTTCAGTCTTT
WI-18369	58	G	0.88	A	0.13	0.22	ATCTGTACACACAATCRAAATGGATAAAGGCCT
WI-18387	57	A	0.63	G	0.38	0.47	TTGGTGACCCCATACRITTTGTGGTCACATGC
WI-18420	38	C	0.19	T	0.81	0.30	GGAAATGGGAAGAAAYAGAGTGAAATTAAG
WI-18420	108	T	0.38	C	0.63	0.47	TCAAAAAAATCAAAAYGCTTATAGCAATGCT
WI-18425	81	A	0.13	C	0.88	0.22	TCCTAGACAGATTTCAMTGCACACAACAACAG
WI-18449	129	C	0.38	T	0.63	0.47	CTCTAAGTGGGACTAYTCTGGATACAGTCAG
WI-18457	120	T	0.94	C	0.06	0.12	ACATTGGGGCCACAGYAAATAGGCTAAAAGG



Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
WI-18462	39	A	0.56	G	0.44	0.49	CAATGGCAGAGGTGARTAGAAACCATCTCAA
WI-18476	60	C	0.56	T	0.44	0.49	GGTGGGGTGCAGGYGGTCACTCCCATCGT
WI-18491	109	G	0.69	A	0.31	0.43	AAATCCAGAAATGACRGGATTACAAGAAAAT
WI-18517	87	C	0.81	T	0.19	0.30	GAATCAGCAGCCTGAYTGTGCACTTGTCCA
WI-18533	59	T	0.44	G	0.56	0.49	TCCCGAGATTTTCTKCTTTATTTTATATT
WI-18533	91	T	0.56	C	0.44	0.49	CATTTTTCATCCTAAAYTTACTGAAGCCATTT
WI-18612	37	A	0.56	G	0.44	0.49	CAAGTTTGGAAATGCRTAATTTGCAAGCAGCA
WI-18640	121	T	0.44	C	0.56	0.49	TGGGGGGTGCAGAGYGTGCTCTTTCAGTG
WI-18668	76	C	0.19	T	0.81	0.30	AAACTAGGCAAAAYAGCAAAAAGTGCAGT
TIGR-A003M18	29	A	0.75	G	0.25	0.38	AGATGAGGTTTTTCCRTGTGTGGCCAGGATGG
TIGR-A003P30	117	C	0.94	G	0.06	0.12	TTTAAAGCAGTGTCSASACTGGCTGCCTGAAG
TIGR-A004S34	156	C	0.25	T	0.75	0.38	CCTCAITCCTATAAAYCTTTAACAAAAACAG
TIGR-A004T44	69	G	0.81	A	0.19	0.30	AACCAAAATGATTGARTATGATAAAGAAATT
TIGR-A004T44	97	A	0.75	C	0.25	0.38	TTTTCATGGCGATTMAAATAGAAAAACCTAT
WI-18673	29	A	0.00	G	1.00	0.00	GTTTTAATTGCAACRACCTTAATTTACAGCA
WI-18680	75	T	0.50	C	0.50	0.50	CTCTAGCATCTGGAAAYGCTCCGTTGTATATT
WI-18694	41	A	0.56	T	0.44	0.49	AGCCAGCTCTGACTTWCTCTCTGTTTCTGTC
WI-18704	99	A	0.56	C	0.44	0.49	TTCTCCGAGGGGTACMCCAGCAGGGCCCTTCA
TIGR-A004V08	60	T	0.88	C	0.13	0.22	ACAGGCATTCTCTTAYGCCCTTTTGTGGGAAG

Fragment	Position	"Ref" Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
TIGR-A004V26	125	A	0.94	G	0.06	0.12	ATTATCTTCACATGARAAAGGTTTCAGTTTAT
TIGR-A004V28	29	A	0.38	G	0.63	0.47	TGTGGTGCGATCTCRGCTCACTGCAACCTC
TIGR-A004X20	25	T	0.31	C	0.69	0.43	TTCTCTTCTGTAGGAYGTCTCCCATGTTACAG
TIGR-A004X30	26	T	0.44	C	0.56	0.49	TAGAGTAGAACCCACACACTCTAGTAATACTT
TIGR-A004Z04	102	T	0.50	G	0.50	0.50	TGGGTATGCAAAACIKTTGCTTTCATGAAAT
TIGR-A004Z19	85	C	0.88	T	0.13	0.22	CATTTTCTTTTCTTTCTTCTCCCGATGACCA
TIGR-A004Z42	89	C	0.88	T	0.13	0.22	GGGAGGTAGGAGACTYGGACCCGGCAGCCCTG
TIGR-A005D17	79	G	0.63	C	0.38	0.47	GGGAAACCCAGCAAGSCTGTCTAGATTCTTC
TIGR-A005D17	81	T	0.56	C	0.44	0.49	GAAACCCAGCAAGGCYGTCTAGATTCTTCTT
TIGR-A005D44	97	G	0.69	T	0.31	0.43	TAAAACTGTTACACTKTTTGTGTTGGCTTTAA
SGC30018	77	C	0.69	T	0.31	0.43	GCACATACTTCAGGCYTGCGGCACCAACCCCA
SGC30036	42	T	0.75	C	0.25	0.38	TAGACAGAGGCATTAYTTTGAAGATCTTTT
SGC30050	103	A	0.31	G	0.69	0.43	CCAGAAAGCTTTACCRCTCTGTCAGTTAAGCT
SGC30055	32	A	0.56	G	0.44	0.49	ATCTTCAGGATAGGTRATAACAGTGTGAAGG
SGC30072	28	C	0.50	T	0.50	0.50	CTTTATTTTGGACAYGTAGCATGTTTAAAC
SGC30076	97	C	0.75	T	0.25	0.38	GGTCACTTTGGGGCCYGGCGTGGGCAGAGCC
SGC30117	96	A	0.50	G	0.50	0.50	GCAGTCACAAATGTACRAAAATGTGACAAGAT
SGC30122	74	A	0.25	G	0.75	0.38	GCAGAACTTAAACACRGAGCATTTATTGTTA
SGC30126	61	T	0.94	C	0.06	0.12	AGTGAATTCAACAGTYAATGCACATGCATAC

17

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC30417	69	T	0.31	C	0.69	0.43	CCTTCCTAAAGCCTCCYAAAGTGCCAGGATTA
SGC30422	24	G	0.44	A	0.56	0.49	CACCTTCTGGAGGCTGRGAAAGTCTAAGATTGA
SGC30515	56	A	0.56	G	0.44	0.49	CAGTACAAAGTCTGTGTRATCCAGGAAAGTGACC
SGC30535	23	T	0.94	A	0.06	0.12	TCCACCTGACCTGCACAWCAACAGCCAGTTAT
SGC30540	52	6A	0.56	G	0.44	0.49	TGTTAAAAACAACACACRTCTGTCACTTGCAGA
SGC30587	74	T	0.94	G	0.06	0.12	AATAGTCTGGCCATTKGACTAACCAAGTTCTA
SGC30593	72	G	0.31	A	0.69	0.43	AGATGTGAGAGACGCRCTCTGTACAGGAGC
SGC30598	70	T	0.81	C	0.19	0.30	GATTTTCTCAGGCCTYTTTTGGGATACCTTTA
SGC30610	99	A	0.50	T	0.50	0.50	CAITTCAGTCCAAGAWAACCTTCCTCCTCAAATT
SGC30612	39	A	0.88	G	0.13	0.22	AAGTTGGGTTTCTTCTGCTGAAATTTCCCATGA
SGC30622	32	T	0.63	C	0.38	0.47	CTGTTACGTCCTTTCCYATTATATTATCTTG
SGC30669	39	A	0.94	G	0.06	0.12	CACAGAGACTGTCTCRGAGACGGGCACAGAA
SGC30678	30	G	0.44	C	0.56	0.49	AATCCCTTGGTGGSGGGGGGGGGTGAGAT
SGC30689	58	G	0.81	A	0.19	0.30	TCATCAGAACCCACACRGTAACCTTGAGTACCT
SGC30719	53	G	0.69	A	0.31	0.43	AGCATCAITGTCACTRGCTAACTCCTCAAAT
SGC30720	85	T	0.63	C	0.38	0.47	CCATCTACAAAAGATYTCCTCATTTGAGGCCTC
SGC30754	66	T	0.81	C	0.19	0.30	CATGTTCTGTTTAAATCTCTTATGTGTTAT
SGC30775	58	A	0.94	G	0.06	0.12	ATTCCAGCAGGTGCCRTTATTTTTCACCTTGGT

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC30813	103	C	0.75	T	0.25	0.38	CCCAGCTTACAGGCTYACAGAAGAAATGAGAC
SGC30827	121	G	0.94	C	0.06	0.12	CTACATAGGGATAAASAGCTCAGTATCTGGA
SGC30890	87	C	0.63	T	0.38	0.47	GTTGTCCAGCCAACAAYGGAGGTGATTTTGGT
SGC30895	72	T	0.31	C	0.69	0.43	AAITTTGTGTCGATGCYCTGTGTCTCCCGTCC
SGC30914	75	T	0.88	C	0.13	0.22	CTATAAGTGCATTTTYYATAATGGGGATTTC
SGC30914	95	T	0.56	G	0.44	0.49	TGGGATTTTCTGTCTKAACTGCCCACTGATT
SGC30938	80	A	0.38	G	0.63	0.47	ATGCAGGAGGGTGGCRAGAGGGGCCGAGATT
SGC30940	103	C	0.94	T	0.06	0.12	AGCTGGCTTTGTAGTYGTTTCAGGCCCATGA
SGC30955	69	A	0.81	G	0.19	0.30	TACTCAAAGTGTGAATRGATTTTATTAGTTGT
SGC30985	75	A	0.75	G	0.25	0.38	AGGACTCTGCATTGTRATTAAAGTTTATTAAAT
SGC31224	47	A	0.88	G	0.13	0.22	AGCATGGCTAAACGRTAAAGATGGGAATCA
SGC31233	85	A	0.63	G	0.38	0.47	AACTTATAACCTCACRCGCTTGTTTCACAAA
SGC31250	79	T	0.75	G	0.25	0.38	TAAGGCCTAAGGAATKAGGGGCAGGGGGCGA
SGC31279	42	G	0.56	A	0.44	0.49	GACCCCTTCGGTGACCCRCAGGCTCCCTGCCAG
SGC31299	57	C	0.31	G	0.69	0.43	TGTCTAAATTTTCCAASACTATGTTTAAATGA
SGC31303	117	C	0.56	T	0.44	0.49	GACTTCAGAGTAATAAYGGTTTATGTCAGTTT
SGC31319	31	C	0.81	T	0.19	0.30	CCCACAAATTTTGATTYGGTGGCTTCATAAGG
SGC31324	45	A	0.81	C	0.19	0.30	GAATAACTGATGTTTCMCAATACCCCGACCCC

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC31372	81	A	0.75	G	0.25	0.38	GACAICTAACCAATTAGRTAGCCTTCAGAAATTG
SGC31485	84	C	0.25	G	0.75	0.38	TTTCTCTAATCTCTSCAAATTGTTTAAAGA
SGC31490	131	C	0.88	T	0.13	0.22	TATATATATGGCTTTTCAATAACCACTAAA
SGC31493	138	G	0.69	A	0.31	0.43	ACTTTGAAATGTAAACRAATGGTACTACAACC
SGC31494	129	T	0.50	C	0.50	0.50	CCTCTGCTGCCATGGYGTGTCCTCTCGGAA
SGC31500	103	C	0.81	T	0.19	0.30	TCTTTTGGACCAACACCTTTTGTCTTTAGAG
SGC31534	159	A	0.75	T	0.25	0.38	AGGAATCTGGGAATTWGCCCCTGGCCTGAAAG
SGC31566	72	T	0.63	G	0.38	0.47	TTTTGTTTATGGATCKGATAAAATCTAGATC
SGC31576	106	G	0.88	C	0.13	0.22	CCAACGATCATATCTSTATGCCTCATTTAT
SGC31596	24	C	0.56	T	0.44	0.49	GTGACGTATGTAGAAYGCTTAGGGTGTCTC
SGC31598	44	C	0.94	T	0.06	0.12	TGCTCTCATCACCAGYTAGAGCTTCTTCCCG
SGC31656	88	G	0.81	A	0.19	0.30	CGACTACCACTGATRAAATACCTGCAAAAGT
SGC31729	128	G	0.88	A	0.13	0.22	CCAATTTAATAAGTGRATATGCTTTCTGAACA
SGC31748	19	A	0.31	C	0.69	0.43	GGAGCTCTGAGGAGCMCAACCAAGGGACGTGT
SGC31767	41	T	0.13	C	0.88	0.22	TATTGAGTTATAATAYACATAAAAAATCCACC
SGC31767	54	A	0.13	G	0.88	0.22	TATACATAAAAAATCCRCCACTGTAAACAGTA
SGC31767	92	T	0.13	C	0.88	0.22	ATGGTTTTTACTCTAYTGTCAAAGCTGGGCA
SGC31772	74	C	0.38	T	0.63	0.47	CGGCACAGACAGAGTYTGGGAGCCATGGGGC

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC31777	118	C	0.88	G	0.13	0.22	GGAGATGCCCCCATGTSTGTGAGACTTAAAA
SGC31788	48	G	0.63	C	0.38	0.47	AAAGACAACAGAGGASAGCAGAGAATAATAT
SGC31914	100	A	0.94	G	0.06	0.12	TTTACATTCAAGGACRGCTTCCAGACAAGCC
SGC31986	61	T	0.44	G	0.56	0.49	TAAAGGGCATAATCTKAAACAAAATTCTTTC
SGC32030	51	A	0.38	G	0.63	0.47	GAACAATATTTTAGGRATTTGAAAATTATTTC
SGC32039	69	G	0.94	C	0.06	0.12	GACTAGTTCAAGCAGSAGGTTAGACCAGTAA
SGC32060	115	T	0.13	G	0.88	0.22	CGGGAGTGCTGATTGKTCGGGTCCCAAGATAA
SGC32109	78	T	0.88	C	0.13	0.22	TGCTATTCCTGCCATYACCGCATCCTTCATG
SGC32119	31	T	0.63	A	0.38	0.47	TGTTTCTTCTTTAAAWATGTTATAAAAATAA
SGC32190	27	C	0.88	T	0.13	0.22	CCAGGCTGGTCTCATYTCAGGCTCATGCGAT
SGC32204	91	T	0.63	C	0.38	0.47	CATTTTTCATCCTAAAYTTACTGAAGCCATTT
SGC32206	40	A	0.25	C	0.75	0.38	TTAAGGGTATAGTTCMAGTGCCATTAAAGTAC
SGC32206	41	A	0.38	G	0.63	0.47	TAAAGGGTATAGTTCARGTGGCATTAAGTACA
SGC32299	108	T	0.94	A	0.06	0.12	ATTAAATCTTTGCCCTTWTATGGTTTTGACAGTT
SGC32391	44	G	0.25	A	0.75	0.38	TTTCAATACTAAACARTGTAAACAATGCAAA
SGC32394	31	T	0.63	G	0.38	0.47	GTTTGTGTTTTTCTCTKTATTTGATGGGATTTA
SGC32407	51	C	0.88	T	0.13	0.22	TGTTCTCCAGTCTTGYAGGTTACATAAGCCA
SGC32411	98	T	0.81	C	0.19	0.30	TTCTCTCAAGTCCCTTTCATCCATACCACCA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC32541	99	A	0.94	G	0.06	0.12	TTATTTTAAATATTCRGGATTAAATTTCTTC
SGC32577	91	T	0.13	G	0.88	0.22	GATGCCAATACTTCGKGCTTCCCAGAGTGCA
SGC32579	24	C	0.44	T	0.56	0.49	CCTAAAGATCTTTTCTCCCAAGTCCTAA
SGC32586	101	C	0.63	G	0.38	0.47	CCTGCTCCGCCCTTCSGCCACCATCCATTCC
SGC32590	86	A	0.63	G	0.38	0.47	CCTGAGGTGATATGGRCCCTTAAGTCCACGAT
SGC32609	72	T	0.94	C	0.06	0.12	ATTCCTAAAATCTATYACACTGAGAGGAAAA
SGC32612	63	G	0.69	T	0.31	0.43	TGAAACAGGGATGCCKTTCTCGGTACTATGT
SGC32620	86	A	0.63	G	0.38	0.47	GATTAGCGTGAGAGGAAAAATGTGAAATGT
SGC32638	97	T	0.63	C	0.38	0.47	TTTCCAGTTGGTAAGYAGCAGGTGCCGAGGG
SGC32641	26	C	0.75	T	0.25	0.38	CGACGCCGCGAGTGTGTGGACCTGTACGTG
SGC32650	83	A	0.69	C	0.31	0.43	CCTTGTTTCAGATTTCMAAATAGTTGTAGCCT
SGC32816	79	C	0.75	T	0.25	0.38	TATATGTGCAGGGCCYGGGGGGTGAAGGGT
SGC32859	78	A	0.13	G	0.88	0.22	TGGAACTGAAACACRGACGCCCTTCTTCCA
SGC32871	39	C	0.94	T	0.06	0.12	TCATCCCAGATTATTYTGAAGTGGAACCCAC
SGC32871	128	C	0.75	T	0.25	0.38	AGACAGTGAGCTGTGTGAGCTGGATTATGTC
SGC32909	26	A	0.75	C	0.25	0.38	GGTAACCCAGTTTGTMACATTATTCAGAACT
SGC32909	95	C	0.44	G	0.56	0.49	GGAGAAAGCAGTGTGSTATAATGTCAACATC
SGC32942	92	G	0.44	C	0.56	0.49	GCGCCGGGCCCTGCCSSGACCCCTGTGTTTCCC



Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC32968	73	T	0.63	G	0.38	0.47	TTTGTGTGTGTGTGTTTAAATTAAGAA
SGC32975	100	T	0.81	A	0.19	0.30	ACCTTCAAAAATTAAWTGTGACTTACGGAAA
SGC32978	58	T	0.88	C	0.13	0.22	CAGCTTGTATTACTTYACAAGTCAGACCTGT
SGC32986	114	A	0.88	G	0.13	0.22	TAGCAGCTTTTAGGGRTTATATCATGAGGTA
SGC32991	56	T	0.88	C	0.13	0.22	ACTGGATAAATAAAAAYGTGGTACATGTACAC
SGC32993	38	C	0.44	T	0.56	0.49	TTTATATAAAACCCTGYAGATGAATATTTTT
SGC33004	53	A	0.63	T	0.38	0.47	ATTTTGGCTAAATTTGTGTAGTCTTACAAAGGC
SGC33092	96	T	0.75	C	0.25	0.38	CATTAAAAATGAACCTYGGAAATAAGAGCATAA
SGC33161	101	G	0.81	T	0.19	0.30	CATTTAAGAAATGAAGKGGAAATGAAGGCAAT
SGC33169	109	C	0.88	T	0.13	0.22	CTCATCTGCTGGTGTGTYTTCCTCAGAGCTTAA
SGC33221	74	A	0.63	G	0.38	0.47	ACTACTCTCCTTCARGACTATTTCATTCTG
SGC33235	82	G	0.94	A	0.06	0.12	CACATAGATCCCAGARTATTAAAGGGCTGG
SGC33289	52	C	0.63	G	0.38	0.47	AGGTCACACTTGTGCASCAGCAAGTATAAACA
SGC33301	95	A	0.50	G	0.50	0.50	ATTAACTGAGATTATRGGAACGCACAGCAA
SGC33302	25	A	0.94	G	0.06	0.12	TTCTGGGCCTGTCAGRAAGTGACATCTTTTA
SGC33319	22	C	0.50	A	0.50	0.50	TCTCCAGGATTCCAGMCTCGTAGCTGATGTG
SGC33355	66	A	0.75	G	0.25	0.38	GCAGTGTCTGGAGACRGTTTTTGATTGTCACA
SGC33366	45	C	0.06	G	0.94	0.12	ATGATTTCAGCATTTASACTTTTAAAAAATTACC

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC33368	69	T	0.75	C	0.25	0.38	GTCTAGGAGTAGAAAYGCACACAGGAATAA
SGC33387	105	T	0.69	C	0.31	0.43	CAGTGTGCCTGAGAYGATGCATGTGGCAGA
SGC33388	67	A	0.31	G	0.69	0.43	CATTGTCCACCGGGCCTTGAGAATAACAATAT
SGC33424	104	C	0.94	T	0.06	0.12	TCTTCTAGGGCCACAYGGAGCAGAACGAGCT
SGC33431	66	A	0.75	G	0.25	0.38	CCAAATAAAATGCACRTAATTTAAAGTTTACA
SGC33436	44	C	0.75	T	0.25	0.38	TAGGTTTGTCTTCCYAGCATATTCAGCTAT
SGC33475	101	C	0.44	G	0.56	0.49	ACTTCTTGAACAAASTGATTACGAAAGTGA
SGC33492	25	A	0.75	G	0.25	0.38	CTGTAAACCGAGCCRCAGTGACCGGGACTT
SGC33497	63	G	0.13	A	0.88	0.22	GCCTCACACAAGCATRATCAATCGCCACGAG
SGC33497	80	G	0.94	A	0.06	0.12	TCAATCGCCACGAGARACTGGATGCCAAAGA
SGC33499	23	A	0.56	G	0.44	0.49	TCCATGTGAACATATRACCTATTCATAAAGT
SGC33533	58	C	0.75	T	0.25	0.38	AATACGAACAGTGCAYGCTGATGGCCTGCAG
SGC33533	102	G	0.88	A	0.13	0.22	TTGGCTCTCTGGACGRTTCATTTCTACATGGC
SGC33565	89	G	0.63	C	0.38	0.47	CAGAAAAGCGCGCTCSGGGTTTTCTGAACCC
SGC33567	96	C	0.75	T	0.25	0.38	CCTAAGTAGTCTCTCYAAAGAGCCATCCCTG
SGC33570	109	C	0.56	A	0.44	0.49	TATAGCCAAGGGACTMGGAAATTTTGGCTGCT
SGC33582	58	T	0.50	C	0.50	0.50	CAGGGCAACATAGGAYTGTGACAGCACCCT
SGC33603	53	C	0.31	T	0.69	0.43	CTAGAGGAGAGATTAYAAATGAACGTAAATAA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC33608	83	A	0.75	G	0.25	0.38	TGGTTCCTCCAGGGGARTTGGCCCCGGAAGCTG
SGC33610	28	G	0.69	A	0.31	0.43	TGGCTTTCAAAATCARTACAGACAGATAAGA
SGC33623	95	C	0.38	A	0.63	0.47	GCCGAGGTCACCTGCTMTACAAAGATTAAAGA
SGC33642	24	C	0.88	G	0.13	0.22	GTGAAGGGACAGAGAGSGTAACACACAGTCCAT
SGC33691	32	C	0.94	G	0.06	0.12	GTTCAGGTTTTTTTTSTGAACAAATGATCCT
SGC33691	119	T	0.63	C	0.38	0.47	CAGTCACCTAAGATAYCGAGTGGCAAAGTCTT
SGC33707	54	G	0.75	A	0.25	0.38	ATGTTAAACAGCAGRAAAATGGAAATTATAGC
SGC33710	59	A	0.19	G	0.81	0.30	AGAAATATCTAGTTGRGTAGAGGAAGGCACT
SGC33712	28	A	0.69	G	0.31	0.43	ATGACACTGCCAACARTCACAGATTTCGATA
SGC33724	45	T	0.69	C	0.31	0.43	GAGTCACAGTTTCATYTGGGAGTCCCTGTGC
SGC33724	52	T	0.44	C	0.56	0.49	AGTTTCATTTGGGAGYCCCTGTGCAGCCCTT
SGC33731	49	G	0.19	A	0.81	0.30	TGTCCCAGTGCCACARTGGTCTAGCCTCATG
SGC33736	62	C	0.44	G	0.56	0.49	TTCAGTTGACAGATTSTCTCCTTACCTAACT
SGC33754	55	C	0.25	T	0.75	0.38	TTGACTCAAGGGCATYGTAAATAGGTTTCCAT
SGC33754	69	A	0.06	G	0.94	0.12	TCGTAATAGGTTTCCRTACTGCAGAGAAGAAGG
SGC33764	71	C	0.25	T	0.75	0.38	GGAAAACAGGAAATCYATCCTTCAAGCATT
SGC33768	41	A	0.13	G	0.88	0.22	AGGCATGAGGAGCTGRTTATGCAGATATACT
SGC33773	38	G	0.38	A	0.63	0.47	ACAACTTGCAAGCACRGGGAGAAAAACCTAGG

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC33835	32	C	0.06	T	0.94	0.12	TTGTCCATTGTTGAYTGTGAATAATTGGCA
SGC33887	118	T	0.63	C	0.38	0.47	TCACATAGGCAGTTGYACACCCAGCTGACAA
SGC33917	72	T	0.50	C	0.50	0.50	GACATGTGTGGCTGYGAGGGAGAGGACCC
SGC33945	21	C	0.75	T	0.25	0.38	TTGCTTAGCCAGCTTYATCAGTGGTGCCCTA
SGC33952	102	C	0.75	G	0.25	0.38	CATAAATTATCAAAATSTGCGCCAGTAATCT
SGC33970	76	C	0.88	G	0.13	0.22	CCCTACTTAGACCCCTSGCACACAAAGGTTGA
SGC33989	31	T	0.88	A	0.13	0.22	AGTCCCTAGGTGTGTWTGAACAATCTGGGT
SGC33991	93	A	0.75	G	0.25	0.38	AAATCACAGTACTGGRATCAGTGAAATTG
SGC34004	90	T	0.75	C	0.25	0.38	AGCAAACCAATAAAAYCATATATCTTGAGGG
SGC34009	46	G	0.50	A	0.50	0.50	TAAGACAGTGTCTACRTGGCCTGAATGTTGG
SGC34014	75	T	0.25	A	0.75	0.38	GGTACCAATATCAATWCAGTTTTTCAAAGCCA
SGC34014	98	T	0.88	C	0.13	0.22	CAAAGCCATTTGCAGYACTCTTCAGATGGGT
SGC34016	44	T	0.94	C	0.06	0.12	AACGGTTTGTAGTTTYGCTTACCCGCAGTGC
SGC34029	53	A	0.56	G	0.44	0.49	CAGATCTGTTTTTCAGRAAGAGGGCCTACTTT
SGC34033	86	C	0.75	T	0.25	0.38	TTTTGACCTATCTCAYCAAGCGAGAGGGAGG
SGC34033	107	G	0.75	A	0.25	0.38	GAGAGGGAGGCAAGCRGAGGGATGGTTTATC
SGC34037	68	A	0.63	G	0.38	0.47	CTGATGGAAGCATCARTGATGGATTTGGCTT
SGC34039	63	T	0.88	C	0.13	0.22	CAGCTTGTGTTGATGYCTACAAAGAAAGTCAG

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC34088	64	T	0.88	C	0.13	0.22	CAAAGCTGAAACCTAAAGAGTGAGCATAGCAA
SGC34119	25	T	0.94	C	0.06	0.12	ATGGAAAGAGTGACAYCCTTGTCCTGTCTG
SGC34142	49	A	0.38	G	0.63	0.47	GAAAACTGATACACCCRGTTACTACTCT
SGC34145	80	G	0.50	A	0.50	0.50	TACTAGGTGCTGGGARTGTGACAGTGAGCAA
SGC34158	26	A	0.44	G	0.56	0.49	AATGACAAAGCCCCAARAGAACAGAGGATCAA
SGC34223	106	T	0.63	C	0.38	0.47	TTTAGCGTAAATACCCYGAATAACCCATAGTT
SGC34226	73	G	0.88	A	0.13	0.22	CAGGCATAAGCAGCCRTGCCCTGACCCACATT
SGC34248	25	T	0.81	C	0.19	0.30	AAAGTAAGCAGCCGGYTGTCCTCGATTGA
SGC34278	33	G	0.81	A	0.19	0.30	GACCTGCTCCTAAAARCTTTCTCCTCCTCCT
SGC34351	51	G	0.25	A	0.75	0.38	CTGTGAACTATGAACRTCTCAGCCTAGAAAGG
SGC34363	28	T	0.13	C	0.88	0.22	CTACCAGAACTCATGYGATAGCGCTTCTTT
SGC34377	78	A	0.31	T	0.69	0.43	GGAAACTTACAATCAWGGTAGAAGGCCAAAAG
SGC34392	56	T	0.50	G	0.50	0.50	GTTTTATATCACTTAKTTATCTCAACAATCT
SGC34411	50	G	0.75	A	0.25	0.38	AGCTCTCAGGACTGGRGCTAGGGTTTAAGGA
SGC34413	59	C	0.75	G	0.25	0.38	AAAGTTCAGTAGAGASAGGTGTTTGAATGT
SGC34485	77	C	0.88	T	0.13	0.22	AAATGATCATTTAACYTCTTTGAACTACAGC
SGC34486	75	C	0.81	G	0.19	0.30	CTTAGAGAAAGTTAASGCACATAGTATTATT
SGC34488	33	G	0.88	T	0.13	0.22	CATGACTACCAACGCKGGCCCTTGCAACCCA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC34489	27	T	0.94	C	0.06	0.12	CTCAAAATCCTAAAYGIGTGCTTCAAAGA
SGC34498	126	A	0.75	C	0.25	0.38	TACACACTGAGCAACMAAAACAAAGGTGTGA
SGC34531	90	A	0.44	G	0.56	0.49	TAAATCGTCTATAGRACCATTTCCCGTCTC
SGC34575	126	C	0.94	G	0.06	0.12	ATTTGATGCAGTTTSGTTAGGGAATTAAGA
SGC34640	97	T	0.67	A	0.33	0.44	GCTGTGGGAACCTCWGGTGCCTTACAACTC
SGC34662	25	G	0.50	A	0.50	0.50	GGAAAAATGGTGGCCTGCCTCTAAAAACCTG
SGC34671	104	A	0.83	G	0.17	0.28	CAGGATGTTCCCTGARGTAATTCAGGAATTCT
SGC34681	82	G	0.08	A	0.92	0.15	TGGGGAGTCTATGTTTGTCTTCTGTTGGCC
SGC34681	93	T	0.92	G	0.08	0.15	TGTTGTGCTTCTGKGCCCTTAAAGAAAC
SGC34724	93	G	0.42	T	0.58	0.49	ATAAAGAGGTTCTCKGCCCTTCCAGCGTTG
SGC34725	33	C	0.83	T	0.17	0.28	TGTAGGCATTAAATGYTATAAAATTTCTGCT
SGC34755	32	C	0.83	T	0.17	0.28	TTAGGCAATGGAAA YAGACTTACTGTATGG
SGC34764	51	C	0.83	G	0.17	0.28	CCCACAAAGGCTCCASATGTTAAACGTTTC
SGC34765	89	C	0.92	A	0.08	0.15	CCCATGAAACCAAGAMCTTGTCTCATGATA
SGC34830	62	C	0.88	T	0.13	0.22	TACTGATTGACAAATGYATATTAGCCAGGTAA
SGC34846	93	C	0.25	T	0.75	0.38	CAGCCATGGCCCCCTGYGCTGATGGAGCTTGT
SGC34858	89	C	0.25	G	0.75	0.38	GTCTGGGGATTCTCTASAGGGGACATATCACA
SGC34906	103	A	0.31	T	0.69	0.43	ATTCAAGCAACAATTWCTTTTATGTTCTTA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC34924	88	T	0.56	C	0.44	0.49	TGCTAGGATTACAGGYGTGAGCCACACACC
SGC34951	37	T	0.94	C	0.06	0.12	AGGGAACTAAGCTCYTCAAAAATAACTGAAA
SGC34953	42	A	0.75	G	0.25	0.38	CACTACGCATGCACARATAAAGTCACATCAA
SGC34961	81	A	0.94	C	0.06	0.12	TGAGCTGGTGGAAAAMGGACTTGGAGACAGC
SGC34964	31	C	0.94	A	0.06	0.12	CATAGTGCCTCTAGTMACCTATGAGGCACTA
SGC34974	68	A	0.94	G	0.06	0.12	TGTAATGCACACCCARTCTGTACTCCACAA
SGC34978	105	T	0.56	G	0.44	0.49	TATATTTTGAAAGTCKCAGGAGAAAATGG
SGC34982	93	G	0.63	T	0.38	0.47	ATGTCACTCTAGGAACKAGTAAACAGGTGTTA
SGC34985	41	T	0.69	C	0.31	0.43	TTCAATTAAATAGTAGYTGAGCGCTGGGGGCT
SGC34985	101	G	0.88	C	0.13	0.22	GTGCTGTGTCCTGCASGCTGTCTCAGGCCAA
SGC34990	63	T	0.56	G	0.44	0.49	TCAATTCTGTGAAAACKAACATGCCTCAAAAA
SGC34994	90	A	0.88	G	0.13	0.22	ATAGTAGGAGTATCTRCCCTGCCCTGCTAGA
SGC35006	45	C	0.56	T	0.44	0.49	ATCCTCCTCAAACTTYAAGGGTGAAAAGCAT
SGC35020	46	G	0.06	A	0.94	0.12	AAATATTAACCCCTCTCTCTCAGGAGTGAC
SGC35053	34	A	0.25	G	0.75	0.38	AGTCATTTATTTACCRGTGTCATGAATTCATTA
SGC35081	57	C	0.56	T	0.44	0.49	TTTTCATGTCACTTAYCGCATGGAAGAACGC
SGC35145	100	T	0.38	C	0.63	0.47	GTCATCCTGACTGACYGTCCCTGCAGTGCCC
SGC35186	90	T	0.63	C	0.38	0.47	TACACAAATGCTATGYAAACAAGTTACTGAA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC35222	70	G	0.88	C	0.13	0.22	GGAAGGAAAGTATGCSGTGTTTAGGGAGAG
SGC35233	21	T	0.50	C	0.50	0.50	CATGGGTGTGACCTCYACAGTACCTCTTCT
SGC35238	159	G	0.19	A	0.81	0.30	GCTCTTCATTCTCACRGGGCGCAACCCCTC
SGC35244	81	A	0.56	G	0.44	0.49	GACCTCCTGTGACCCRTGAATGTGCTCCAA
SGC35245	166	C	0.88	T	0.13	0.22	TAATACGTACTTTAGYTGGAATTATTCTATG
SGC35252	39	T	0.31	C	0.69	0.43	TCAGGAACACCCCCAYGACATTGCATTGGG
SGC35267	134	T	0.63	C	0.38	0.47	TTATCCAACTCTCGAYTTTTCCTTGGTCTCC
SGC35276	39	T	0.88	C	0.13	0.22	ATCTGTATTGACTAA YACACCAGTCCACACT
SGC35282	45	C	0.44	G	0.56	0.49	TATCCCTTTTCTCTCTSCAAATGTTTCTCCTC
SGC35282	157	A	0.31	G	0.69	0.43	TTTTTCTTTTCTCARGTGTTACCTACTAAG
SGC35282	173	A	0.50	G	0.50	0.50	GTGTTACCTACTAAGRGATGCCTGGAGTAAG
SGC35285	63	T	0.63	C	0.38	0.47	TCATGTGAAAACTACYCCAGTGGCTGACTGA
SGC35326	34	G	0.81	A	0.19	0.30	TCAGGCTGACGGGGARGAACCACCTGCACCAC
SGC35336	36	T	0.88	C	0.13	0.22	CCTTAGGGCTACAGYCTCTTGTCTGGACC
SGC35345	137	G	0.50	C	0.50	0.50	CAGCGTCCCCCACCSCGTCGTGGTGTAGTC
SGC35346	133	A	0.75	G	0.25	0.38	TGACTGCATGAATGCGTGTGCGTGCAAGCAT
SGC35357	123	T	0.94	G	0.06	0.12	TTGTATTTTGTATATKCGCCTGAAGATCATC
SGC35364	21	A	0.56	G	0.44	0.49	CATCCTGATGCCCCARGTTATCCACAGCCTC



Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC33364	85	A	0.88	G	0.13	0.22	CATTTTCCTGTAAATTCCTCCAACTGATCC
SGC33370	162	T	0.25	C	0.75	0.38	GAACAGCCAGAGATYTTACCGTGGTCTTAC
SGC33384	58	T	0.81	C	0.19	0.30	GTAGTTCTAGAACTTYAGAAAGCTCCATCTTT
SGC33405	114	T	0.75	G	0.25	0.38	TCCCCGACAGCAAAKGGTTTCTCTCTGAGG
SGC33413	62	T	0.88	G	0.13	0.22	TTGTATATAAGATAAKTCATACTACTGGAGAAA
SGC33416	143	C	0.88	G	0.13	0.22	AAAATTGCAAAAAGAMAAAGTATGACTTTTAT
SGC33416	164	C	0.75	A	0.25	0.38	AACCTCAACCACATCYTATCCTCCACCCCAC
SGC33419	25	T	0.31	C	0.69	0.43	CCCTTCCTGGAGACTRAAACCTGGTGTCTCAGG
SGC33432	147	G	0.19	A	0.81	0.30	TATGTTAATTTGCTCTRATACAAAATTCCTAA
SGC33438	99	A	0.88	G	0.13	0.22	GGGAGGGGGCGTTTCRCCTTCCTCTCTCTTG
SGC33461	82	G	0.88	A	0.13	0.22	GGGAGGGGGCGTTTCRCCTTCCTCTCTCTTG
SGC33464	128	T	0.69	G	0.31	0.43	CCGCAAGATGGGGCKGGGCAATGCGCAGGAG
SGC33477	179	G	0.50	T	0.50	0.50	ACAAAGATGGAATTTAKCAAACCCCTAGCCTTG
SGC33498	173	G	0.75	A	0.25	0.38	ACCACAAATCTGAACRTGCCTCTCCCTTGCT
SGC33499	68	C	0.50	T	0.50	0.50	CTTAGGGCATCGCTCYTCTCCTCACGCCACAAA
SGC33499	76	G	0.75	A	0.25	0.38	ATCGCTCCTCCTCACRCCACAAAATCTGGTGC
SGC33527	82	T	0.38	C	0.63	0.47	GGGAAAGTCTGGTCCYACATCTGCCCGCCCT
SGC33529	54	A	0.63	G	0.38	0.47	GGCCCTGAGCGTCTCTRCCCGGAATTCACGAG

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC35531	57	G	0.06	C	0.94	0.12	CACAACCACTTGACSAATGCTTGCCAAGCT
SGC35537	187	T	0.81	C	0.19	0.30	GCAGTCTGGTCCATGYTGGTCTCATACCTCA
SGC35543	78	C	0.81	T	0.19	0.30	CCTCCAGACCGCAGGYTCCCCCAGCCTCAGG
SGC35548	61	A	0.06	G	0.94	0.12	GGTGTGACACACCCARTTTTGTGAGTGTACTGT
SGC35566	70	G	0.88	A	0.13	0.22	TGCAACCAAAACAGCCRTCATCAAAACCCCTCA
SGC35566	78	A	0.44	G	0.56	0.49	AACAGCCGTTCATCAARCCCTCACTAAAAAGT
SGC35569	28	A	0.56	G	0.44	0.49	TCCCAGGCCCCAGGCRCTCTTTCCTGCCCTGC
SGC35569	99	G	0.63	C	0.38	0.47	TCAGCTACTTCTCTCCTSCACTTTTGAAAGACCC
SGC35579	29	C	0.31	T	0.69	0.43	AATTAGCCCTAAATGYGGGTAATAATTTTCC
SGC35580	64	T	0.81	C	0.19	0.30	AATGCAATTGAGCTGYCCAGGCTCTGTCTC
SGC35587	118	A	0.25	C	0.75	0.38	ACATTCAAAAGAAAMGTTGCGAAAATTGCG
SGC35587	148	C	0.63	T	0.38	0.47	GAAATCTGTTGTGCAYGCTCAAAATGAAAACG
SGC35590	191	A	0.56	C	0.44	0.49	TGCAGCTTAAAGAGCMCAGGTTCCAGTACTG
SGC35597	69	T	0.19	C	0.81	0.30	CACCTTCCAAGGCCCYATCCATTAGTTTCCA
SGC35598	24	A	0.00	G	1.00	0.00	TGTCTTCTCTCCACRTGCACAGCTTCCCTGA
SGC35601	113	T	0.56	C	0.44	0.49	CTCCCCATGTGCCTGYGCCAAGAGACAGACA
SGC35615	52	G	0.63	A	0.38	0.47	TATTGTACCAGAACTRTTATTTACACCCCAT
SGC35626	106	A	0.69	T	0.31	0.43	CATGTGGTTTTAAAAAWATCCATAAGGGAAGG

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC35638	20	C	0.19	G	0.81	0.30	ATGAGGCCCATCTCTSGCTCTGTGTGAAG
SGC35645	122	T	0.81	C	0.19	0.30	GCGATGACACACACACACACACACATTTA
SGC35655	101	C	0.25	T	0.75	0.38	ATTTTCTCTGTTCCAYGAAGAGGACTTTTG
SGC35659	150	G	0.19	A	0.81	0.30	TCTTTTCTCCCAAGCRAAACCAAAATGCGCCC
SGC35665	34	A	0.88	G	0.13	0.22	TCAGGAGTCATTAGCRGTGATGATTTTGGGA
SGC35665	89	A	0.31	G	0.69	0.43	TTCCACGTTAGCCARTTGTTCCTTGATGAAT
SGC35671	111	T	0.25	C	0.75	0.38	GGTTTACTTTCAGAAAYGAAGAACTTATTCAG
SGC35678	34	C	0.88	T	0.13	0.22	TGACTCTGCTTCTCTGYACTGACCCAGAGCCT
SGC35687	70	T	0.56	A	0.44	0.49	ATCTCTAAATAAGATWACATTCTGGGGTACT
SGC35825	57	C	0.63	T	0.38	0.47	TAGTAATAAATTACAYGAGATATTCACACTT
SGC35842	98	A	0.94	G	0.06	0.12	ATCCATTATTACAGRAAAATGTGGAAAAGAT
SGC35914	59	T	0.88	C	0.13	0.22	TTATTATGAGCCCCYAGGACCAACATGT
SGC35927	71	C	0.06	T	0.94	0.12	TTCAGTATCATTAATGYTGTAGATTTTCAGATG
SGC35928	25	T	0.50	C	0.50	0.50	TTATCAAAATGGTTAYAGTTTTCAAATTAATA
SGC35946	45	A	0.25	G	0.75	0.38	AATTTTCTCAACTTTCATTTTAAATGTAT
SGC35965	25	A	0.50	C	0.50	0.50	GACATACATATCTCAMGTAGAAATTAGCTATA
SGC35978	36	A	0.88	G	0.13	0.22	ACTTTTTTATAAAGARTAAAGTTGACTGAAAA
SGC35978	45	C	0.50	T	0.50	0.50	TAAAGAATAAGTTGAYTTGAAAAAGCAGTTTAA

Fragment	Position	"Ref Allele"	"Freq (P)"	"Alt Allele"	"Freq (Q)"	"H"	"Sequence Tag"
SGC36020	26	T	0.56	C	0.44	0.49	ACAAGACAATTGCAATYTAACATTGTTATAAA
SGC36047	31	T	0.63	A	0.38	0.47	AGACGGACATAAAAAAATATACACAAAAAAC

### Analysis of Polymorphisms

#### A. Preparation of Samples

Polymorphisms are detected in a target nucleic acid from an individual being analyzed. For assay of genomic DNA, virtually any biological sample (other than pure red blood cells) is suitable. For example, convenient tissue samples include whole blood, semen, saliva, tears, urine, fecal material, sweat, buccal, skin and hair. For assay of cDNA or mRNA, the tissue sample must be obtained from an organ in which the target nucleic acid is expressed. For example, if the target nucleic acid is a cytochrome P450, the liver is a suitable source.

Many of the methods described below require amplification of DNA from target samples. This can be accomplished by e.g., PCR. See generally *PCR Technology: Principles and Applications for DNA Amplification* (ed. H.A. Erlich, Freeman Press, NY, NY, 1992); *PCR Protocols: A Guide to Methods and Applications* (eds. Innis, et al., Academic Press, San Diego, CA, 1990); Mattila et al., *Nucleic Acids Res.* 19, 4967 (1991); Eckert et al., *PCR Methods and Applications* 1, 17 (1991); PCR (eds. McPherson et al., IRL Press, Oxford); and U.S. Patent 4,683,202 (each of which is incorporated by reference for all purposes).

Other suitable amplification methods include the ligase chain reaction (LCR) (see Wu and Wallace, *Genomics* 4, 560 (1989), Landegren et al., *Science* 241, 1077 (1988), transcription amplification (Kwoh et al., *Proc. Natl. Acad. Sci. USA* 86, 1173 (1989)), and self-sustained sequence replication (Guatelli et al., *Proc. Nat. Acad. Sci. USA*, 87, 1874 (1990)) and nucleic acid based sequence amplification (NASBA). The latter two amplification methods involve isothermal reactions based on isothermal transcription, which produce both single stranded RNA (ssRNA) and double stranded DNA (dsDNA) as the amplification products in a ratio of about 30 or 100 to 1, respectively.

## B. Detection of Polymorphisms in Target DNA

There are two distinct types of analysis depending whether a polymorphism in question has already been characterized. The first type of analysis is sometimes referred to as de novo characterization. This analysis compares target sequences in different individuals to identify points of variation, i.e., polymorphic sites. By analyzing a groups of individuals representing the greatest ethnic diversity among humans and greatest breed and species variety in plants and animals, patterns characteristic of the most common alleles/haplotypes of the locus can be identified, and the frequencies of such populations in the population determined. Additional allelic frequencies can be determined for subpopulations characterized by criteria such as geography, race, or gender. The de novo identification of the polymorphisms of the invention is described in the Examples section. The second type of analysis is determining which form(s) of a characterized polymorphism are present in individuals under test. There are a variety of suitable procedures, which are discussed in turn.

### 1. Allele-Specific Probes

The design and use of allele-specific probes for analyzing polymorphisms is described by e.g., Saiki et al., Nature 324, 163-166 (1986); Dattagupta, EP 235,726, Saiki, WO 89/11548. Allele-specific probes can be designed that hybridize to a segment of target DNA from one individual but do not hybridize to the corresponding segment from another individual due to the presence of different polymorphic forms in the respective segments from the two individuals. Hybridization conditions should be sufficiently stringent that there is a significant difference in hybridization intensity between alleles, and preferably an essentially binary response, whereby a probe hybridizes to only one of the alleles. Some probes are designed to hybridize to a segment of target DNA such that the polymorphic site aligns with a central position (e.g., in a 15 mer at the 7 position; in a 16 mer, at either the 8 or 9 position) of the probe. This design

of probe achieves good discrimination in hybridization between different allelic forms.

Allele-specific probes are often used in pairs, one member of a pair showing a perfect match to a reference form of a target sequence and the other member showing a perfect match to a variant form. Several pairs of probes can then be immobilized on the same support for simultaneous analysis of multiple polymorphisms within the same target sequence.

## 2. Tiling Arrays

The polymorphisms can also be identified by hybridization to nucleic acid arrays, some example of which are described by WO 95/11995 (incorporated by reference in its entirety for all purposes). One form of such arrays is described in the Examples section in connection with de novo identification of polymorphisms. The same array or a different array can be used for analysis of characterized polymorphisms. WO 95/11995 also describes subarrays that are optimized for detection of a variant forms of a precharacterized polymorphism. Such a subarray contains probes designed to be complementary to a second reference sequence, which is an allelic variant of the first reference sequence. The second group of probes is designed by the same principles as described in the Examples except that the probes exhibit complementarity to the second reference sequence. The inclusion of a second group (or further groups) can be particularly useful for analyzing short subsequences of the primary reference sequence in which multiple mutations are expected to occur within a short distance commensurate with the length of the probes (i.e., two or more mutations within 9 to 21 bases).

## 3. Allele-Specific Primers

An allele-specific primer hybridizes to a site on target DNA overlapping a polymorphism and only primes amplification of an allelic form to which the primer exhibits perfect complementarity. See Gibbs, *Nucleic Acid Res.* 17, 2427-2448 (1989). This primer is used in conjunction with a second primer which hybridizes at a distal site. Amplification proceeds from the two primers leading to a

detectable product signifying the particular allelic form is present. A control is usually performed with a second pair of primers, one of which shows a single base mismatch at the polymorphic site and the other of which exhibits perfect complementarily to a distal site. The single-base mismatch prevents amplification and no detectable product is formed. The method works best when the mismatch is included in the 3'-most position of the oligonucleotide aligned with the polymorphism because this position is most destabilizing to elongation from the primer. See, e.g., WO 93/22456.

#### 4. Direct-Sequencing

The direct analysis of the sequence of polymorphisms of the present invention can be accomplished using either the dideoxy chain termination method or the Maxam Gilbert method (see Sambrook et al., *Molecular Cloning, A Laboratory Manual* (2nd Ed., CSHP, New York 1989); Zyskind et al., *Recombinant DNA Laboratory Manual*, (Acad. Press, 1988)).

#### 5. Denaturing Gradient Gel Electrophoresis

Amplification products generated using the polymerase chain reaction can be analyzed by the use of denaturing gradient gel electrophoresis. Different alleles can be identified based on the different sequence-dependent melting properties and electrophoretic migration of DNA in solution. Erlich, ed., *PCR Technology, Principles and Applications for DNA Amplification*, (W.H. Freeman and Co, New York, 1992), Chapter 7.

#### 6. Single-Strand Conformation Polymorphism Analysis

Alleles of target sequences can be differentiated using single-strand conformation polymorphism analysis, which identifies base differences by alteration in electrophoretic migration of single stranded PCR products, as described in Orita et al., *Proc. Nat. Acad. Sci.* 86, 2766-2770 (1989). Amplified PCR products can be generated as described above, and heated or otherwise denatured, to form single stranded amplification products. Single-stranded nucleic acids may refold or form secondary structures which are partially dependent on the base sequence. The different electrophoretic mobilities of single-stranded amplification products can be



related to base-sequence difference between alleles of target sequences.

### III. Methods of Use

After determining polymorphic form(s) present in an individual at one or more polymorphic sites, this information can be used in a number of methods.

#### A. Forensics

Determination of which polymorphic forms occupy a set of polymorphic sites in an individual identifies a set of polymorphic forms that distinguishes the individual. See generally National Research Council, *The Evaluation of Forensic DNA Evidence* (Eds. Pollard et al., National Academy Press, DC, 1996). The more sites that are analyzed the lower the probability that the set of polymorphic forms in one individual is the same as that in an unrelated individual. Preferably, if multiple sites are analyzed, the sites are unlinked. Thus, polymorphisms of the invention are often used in conjunction with polymorphisms in distal genes. Preferred polymorphisms for use in forensics are diallelic because the population frequencies of two polymorphic forms can usually be determined with greater accuracy than those of multiple polymorphic forms at multi-allelic loci.

The capacity to identify a distinguishing or unique set of forensic markers in an individual is useful for forensic analysis. For example, one can determine whether a blood sample from a suspect matches a blood or other tissue sample from a crime scene by determining whether the set of polymorphic forms occupying selected polymorphic sites is the same in the suspect and the sample. If the set of polymorphic markers does not match between a suspect and a sample, it can be concluded (barring experimental error) that the suspect was not the source of the sample. If the set of markers does match, one can conclude that the DNA from the suspect is consistent with that found at the crime scene. If frequencies of the polymorphic forms at the loci tested have been determined (e.g., by analysis of a suitable population of individuals), one can perform a statistical analysis to

determine the probability that a match of suspect and crime scene sample would occur by chance.

$p(ID)$  is the probability that two random individuals have the same polymorphic or allelic form at a given polymorphic site. In diallelic loci, four genotypes are possible: AA, AB, BA, and BB. If alleles A and B occur in a haploid genome of the organism with frequencies  $x$  and  $y$ , the probability of each genotype in a diploid organism are (see WO 95/12607):

10 Homozygote:  $p(AA) = x^2$   
 Homozygote:  $p(BB) = y^2 = (1-x)^2$   
 Single Heterozygote:  $p(AB) = p(BA) = xy = x(1-x)$   
 Both Heterozygotes:  $p(AB+BA) = 2xy = 2x(1-x)$

The probability of identity at one locus (i.e., the probability that two individuals, picked at random from a population will have identical polymorphic forms at a given locus) is given by the equation:

$$p(ID) = (x^2)^2 + (2xy)^2 + (y^2)^2.$$

These calculations can be extended for any number of polymorphic forms at a given locus. For example, the probability of identity  $p(ID)$  for a 3-allele system where the alleles have the frequencies in the population of  $x$ ,  $y$  and  $z$ , respectively, is equal to the sum of the squares of the genotype frequencies:

$$25 \quad p(ID) = x^4 + (2xy)^2 + (2yz)^2 + (2xz)^2 + z^4 + y^4$$

In a locus of  $n$  alleles, the appropriate binomial expansion is used to calculate  $p(ID)$  and  $p(exc)$ .

The cumulative probability of identity ( $\text{cum } p(ID)$ ) for each of multiple unlinked loci is determined by multiplying the probabilities provided by each locus.

$$\text{cum } p(ID) = p(ID1)p(ID2)p(ID3) \dots p(IDn)$$

The cumulative probability of non-identity for  $n$  loci (i.e., the probability that two random individuals will be different at 1 or more loci) is given by the equation:

$$35 \quad \text{cum } p(\text{nonID}) = 1 - \text{cum } p(ID).$$

If several polymorphic loci are tested, the cumulative probability of non-identity for random individuals becomes very high (e.g., one billion to one). Such probabilities can

be taken into account together with other evidence in determining the guilt or innocence of the suspect.

#### B. Paternity Testing

The object of paternity testing is usually to  
 5 determine whether a male is the father of a child. In most cases, the mother of the child is known and thus, the mother's contribution to the child's genotype can be traced. Paternity testing investigates whether the part of the child's genotype not attributable to the mother is consistent with that of the  
 10 putative father. Paternity testing can be performed by analyzing sets of polymorphisms in the putative father and the child.

If the set of polymorphisms in the child attributable to the father does not match the putative father, it can be  
 15 concluded, barring experimental error, that the putative father is not the real father. If the set of polymorphisms in the child attributable to the father does match the set of polymorphisms of the putative father, a statistical calculation can be performed to determine the probability of  
 20 coincidental match.

The probability of parentage exclusion (representing the probability that a random male will have a polymorphic form at a given polymorphic site that makes him incompatible as the father) is given by the equation (see WO 95/12607):

$$25 \quad p(\text{exc}) = xy(1-xy)$$

where x and y are the population frequencies of alleles A and B of a diallelic polymorphic site.

(At a triallelic site  $p(\text{exc}) = xy(1-xy) + yz(1-yz) + xz(1-xz) + 3xyz(1-xyz)$ ), where x, y and z are the respective  
 30 population frequencies of alleles A, B and C).

The probability of non-exclusion is

$$p(\text{non-exc}) = 1-p(\text{exc})$$

The cumulative probability of non-exclusion (representing the value obtained when n loci are used) is  
 35 thus:

$$\text{cum } p(\text{non-exc}) = p(\text{non-exc1})p(\text{non-exc2})p(\text{non-exc3}) \dots p(\text{non-excn})$$

The cumulative probability of exclusion for  $n$  loci (representing the probability that a random male will be excluded)

$$\text{cum } p(\text{exc}) = 1 - \text{cum } p(\text{non-exc}).$$

- 5           If several polymorphic loci are included in the analysis, the cumulative probability of exclusion of a random male is very high. This probability can be taken into account in assessing the liability of a putative father whose polymorphic marker set matches the child's polymorphic marker set attributable to his/her father.

10           C. Correlation of Polymorphisms with Phenotypic Traits

- The polymorphisms of the invention may contribute to the phenotype of an organism in different ways. Some polymorphisms occur within a protein coding sequence and contribute to phenotype by affecting protein structure. The effect may be neutral, beneficial or detrimental, or both beneficial and detrimental, depending on the circumstances. For example, a heterozygous sickle cell mutation confers resistance to malaria, but a homozygous sickle cell mutation is usually lethal. Other polymorphisms occur in noncoding regions but may exert phenotypic effects indirectly via influence on replication, transcription, and translation. A single polymorphism may affect more than one phenotypic trait. Likewise, a single phenotypic trait may be affected by polymorphisms in different genes. Further, some polymorphisms predispose an individual to a distinct mutation that is causally related to a certain phenotype.

- Phenotypic traits include diseases that have known but hitherto unmapped genetic components (e.g., agammaglobulinemia, diabetes insipidus, Lesch-Nyhan syndrome, muscular dystrophy, Wiskott-Aldrich syndrome, Fabry's disease, familial hypercholesterolemia, polycystic kidney disease, hereditary spherocytosis, von Willebrand's disease, tuberous sclerosis, hereditary hemorrhagic telangiectasia, familial colonic polyposis, Ehlers-Danlos syndrome, osteogenesis imperfecta, and acute intermittent porphyria). Phenotypic traits also include symptoms of, or susceptibility to, multifactorial diseases of which a component is or may be

genetic, such as autoimmune diseases, inflammation, cancer, diseases of the nervous system, and infection by pathogenic microorganisms. Some examples of autoimmune diseases include rheumatoid arthritis, multiple sclerosis, diabetes (insulin-  
5 dependent and non-independent), systemic lupus erythematosus and Graves disease. Some examples of cancers include cancers of the bladder, brain, breast, colon, esophagus, kidney, leukemia, liver, lung, oral cavity, ovary, pancreas, prostate, skin, stomach and uterus. Phenotypic traits also include  
10 characteristics such as longevity, appearance (e.g., baldness, obesity), strength, speed, endurance, fertility, and susceptibility or receptivity to particular drugs or therapeutic treatments.

Correlation is performed for a population of  
15 individuals who have been tested for the presence or absence of a phenotypic trait of interest and for polymorphic markers sets. To perform such analysis, the presence or absence of a set of polymorphisms (i.e. a polymorphic set) is determined for a set of the individuals, some of whom exhibit a  
20 particular trait, and some of which exhibit lack of the trait. The alleles of each polymorphism of the set are then reviewed to determine whether the presence or absence of a particular allele is associated with the trait of interest. Correlation can be performed by standard statistical methods such as a  $\chi$ -  
25 squared test and statistically significant correlations between polymorphic form(s) and phenotypic characteristics are noted. For example, it might be found that the presence of allele A1 at polymorphism A correlates with heart disease. As  
30 a further example, it might be found that the combined presence of allele A1 at polymorphism A and allele B1 at polymorphism B correlates with increased milk production of a farm animal.

Such correlations can be exploited in several ways. In the case of a strong correlation between a set of one or  
35 more polymorphic forms and a disease for which treatment is available, detection of the polymorphic form set in a human or animal patient may justify immediate administration of treatment, or at least the institution of regular monitoring

of the patient. Detection of a polymorphic form correlated with serious disease in a couple contemplating a family may also be valuable to the couple in their reproductive decisions. For example, the female partner might elect to  
 5 undergo in vitro fertilization to avoid the possibility of transmitting such a polymorphism from her husband to her offspring. In the case of a weaker, but still statistically significant correlation between a polymorphic set and human disease, immediate therapeutic intervention or monitoring may  
 10 not be justified. Nevertheless, the patient can be motivated to begin simple life-style changes (e.g., diet, exercise) that can be accomplished at little cost to the patient but confer potential benefits in reducing the risk of conditions to which the patient may have increased susceptibility by virtue of  
 15 variant alleles. Identification of a polymorphic set in a patient correlated with enhanced receptiveness to one of several treatment regimes for a disease indicates that this treatment regime should be followed.

For animals and plants, correlations between  
 20 characteristics and phenotype are useful for breeding for desired characteristics. For example, Beitz et al., US 5,292,639 discuss use of bovine mitochondrial polymorphisms in a breeding program to improve milk production in cows. To evaluate the effect of mtDNA D-loop sequence polymorphism on  
 25 milk production, each cow was assigned a value of 1 if variant or 0 if wildtype with respect to a prototypical mitochondrial DNA sequence at each of 17 locations considered. Each production trait was analyzed individually with the following animal model:

30 
$$Y_{ijkpn} = \mu + YS_i + P_j + X_k + \beta_1 + \dots + \beta_{17} + PE_n + a_n + e_p$$
 where  $Y_{ijkpn}$  is the milk, fat, fat percentage, SNF, SNF percentage, energy concentration, or lactation energy record;  $\mu$  is an overall mean;  $YS_i$  is the effect common to all cows calving in year-season;  $X_k$  is the effect common to cows in  
 35 either the high or average selection line;  $\beta_1$  to  $\beta_{17}$  are the binomial regressions of production record on mtDNA D-loop sequence polymorphisms;  $PE_n$  is permanent environmental effect common to all records of cow  $n$ ;  $a_n$  is effect of animal  $n$  and

is composed of the additive genetic contribution of sire and dam breeding values and a Mendelian sampling effect; and  $e_p$  is a random residual. It was found that eleven of seventeen polymorphisms tested influenced at least one production trait.

5 Bovines having the best polymorphic forms for milk production at these eleven loci are used as parents for breeding the next generation of the herd.

D. Genetic Mapping of Phenotypic Traits

The previous section concerns identifying correlations between phenotypic traits and polymorphisms that directly or indirectly contribute to those traits. The present section describes identification of a physical linkage between a genetic locus associated with a trait of interest and polymorphic markers that are not associated with the trait, but are in physical proximity with the genetic locus responsible for the trait and co-segregate with it. Such analysis is useful for mapping a genetic locus associated with a phenotypic trait to a chromosomal position, and thereby cloning gene(s) responsible for the trait. See Lander et al., *Proc. Natl. Acad. Sci. (USA)* 83, 7353-7357 (1986); Lander et al., *Proc. Natl. Acad. Sci. (USA)* 84, 2363-2367 (1987); Donis-Keller et al., *Cell* 51, 319-337 (1987); Lander et al., *Genetics* 121, 185-199 (1989)). Genes localized by linkage can be cloned by a process known as directional cloning. See Wainwright, *Med. J. Australia* 159, 170-174 (1993); Collins, *Nature Genetics* 1, 3-6 (1992) (each of which is incorporated by reference in its entirety for all purposes).

10  
15  
20  
25

Linkage studies are typically performed on members of a family. Available members of the family are characterized for the presence or absence of a phenotypic trait and for a set of polymorphic markers. The distribution of polymorphic markers in an informative meiosis is then analyzed to determine which polymorphic markers co-segregate with a phenotypic trait. See, e.g., Kerem et al., *Science* 245, 1073-1080 (1989); Monaco et al., *Nature* 316, 842 (1985); Yamoka et al., *Neurology* 40, 222-226 (1990); Rossiter et al., *FASEB Journal* 5, 21-27 (1991).

30  
35

Linkage is analyzed by calculation of LOD (log of the odds) values. A lod value is the relative likelihood of obtaining observed segregation data for a marker and a genetic locus when the two are located at a recombination fraction  $\theta$ , versus the situation in which the two are not linked, and thus segregating independently (Thompson & Thompson, *Genetics in Medicine* (5th ed, W.B. Saunders Company, Philadelphia, 1991); Strachan, "Mapping the human genome" in *The Human Genome* (BIOS Scientific Publishers Ltd, Oxford), Chapter 4). A series of likelihood ratios are calculated at various recombination fractions ( $\theta$ ), ranging from  $\theta = 0.0$  (coincident loci) to  $\theta = 0.50$  (unlinked). Thus, the likelihood at a given value of  $\theta$  is: probability of data if loci linked at  $\theta$  to probability of data if loci unlinked. The computed likelihoods are usually expressed as the  $\log_{10}$  of this ratio (i.e., a lod score). For example, a lod score of 3 indicates 1000:1 odds against an apparent observed linkage being a coincidence. The use of logarithms allows data collected from different families to be combined by simple addition. Computer programs are available for the calculation of lod scores for differing values of  $\theta$  (e.g., LIPED, MLINK (Lathrop, *Proc. Nat. Acad. Sci. (USA)* 81, 3443-3446 (1984)). For any particular lod score, a recombination fraction may be determined from mathematical tables. See Smith et al., *Mathematical tables for research workers in human genetics* (Churchill, London, 1961); Smith, *Ann. Hum. Genet.* 32, 127-150 (1968). The value of  $\theta$  at which the lod score is the highest is considered to be the best estimate of the recombination fraction.

Positive lod score values suggest that the two loci are linked, whereas negative values suggest that linkage is less likely (at that value of  $\theta$ ) than the possibility that the two loci are unlinked. By convention, a combined lod score of +3 or greater (equivalent to greater than 1000:1 odds in favor of linkage) is considered definitive evidence that two loci are linked. Similarly, by convention, a negative lod score of -2 or less is taken as definitive evidence against linkage of the two loci being compared. Negative linkage data are useful in excluding a chromosome or a segment thereof from



consideration. The search focuses on the remaining non-excluded chromosomal locations.

#### IV. Modified Polypeptides and Gene Sequences

The invention further provides variant forms of nucleic acids and corresponding proteins. The nucleic acids comprise one of the sequences described in Table 1, column 8, in which the polymorphic position is occupied by one of the alternative bases for that position. Some nucleic acid encode full-length variant forms of proteins. Similarly, variant proteins have the prototypical amino acid sequences of encoded by nucleic acid sequence shown in Table 1, column 8, (read so as to be in-frame with the full-length coding sequence of which it is a component) except at an amino acid encoded by a codon including one of the polymorphic positions shown in the Table. That position is occupied by the amino acid coded by the corresponding codon in any of the alternative forms shown in the Table.

Variant genes can be expressed in an expression vector in which a variant gene is operably linked to a native or other promoter. Usually, the promoter is a eukaryotic promoter for expression in a mammalian cell. The transcription regulation sequences typically include a heterologous promoter and optionally an enhancer which is recognized by the host. The selection of an appropriate promoter, for example trp, lac, phage promoters, glycolytic enzyme promoters and tRNA promoters, depends on the host selected. Commercially available expression vectors can be used. Vectors can include host-recognized replication systems, amplifiable genes, selectable markers, host sequences useful for insertion into the host genome, and the like.

The means of introducing the expression construct into a host cell varies depending upon the particular construction and the target host. Suitable means include fusion, conjugation, transfection, transduction, electroporation or injection, as described in Sambrook, *supra*. A wide variety of host cells can be employed for expression of the variant gene, both prokaryotic and eukaryotic. Suitable host cells include

bacteria such as *E. coli*, yeast, filamentous fungi, insect cells, mammalian cells, typically immortalized, e.g., mouse, CHO, human and monkey cell lines and derivatives thereof. Preferred host cells are able to process the variant gene product to produce an appropriate mature polypeptide. Processing includes glycosylation, ubiquitination, disulfide bond formation, general post-translational modification, and the like.

The protein may be isolated by conventional means of protein biochemistry and purification to obtain a substantially pure product, i.e., 80, 95 or 99% free of cell component contaminants, as described in Jacoby, *Methods in Enzymology* Volume 104, Academic Press, New York (1984); Scopes, *Protein Purification, Principles and Practice*, 2nd Edition, Springer-Verlag, New York (1987); and Deutscher (ed), *Guide to Protein Purification, Methods in Enzymology*, Vol. 182 (1990). If the protein is secreted, it can be isolated from the supernatant in which the host cell is grown. If not secreted, the protein can be isolated from a lysate of the host cells.

The invention further provides transgenic nonhuman animals capable of expressing an exogenous variant gene and/or having one or both alleles of an endogenous variant gene inactivated. Expression of an exogenous variant gene is usually achieved by operably linking the gene to a promoter and optionally an enhancer, and microinjecting the construct into a zygote. See Hogan et al., "Manipulating the Mouse Embryo, A Laboratory Manual," Cold Spring Harbor Laboratory. Inactivation of endogenous variant genes can be achieved by forming a transgene in which a cloned variant gene is inactivated by insertion of a positive selection marker. See Capecchi, *Science* 244, 1288-1292 (1989). The transgene is then introduced into an embryonic stem cell, where it undergoes homologous recombination with an endogenous variant gene. Mice and other rodents are preferred animals. Such animals provide useful drug screening systems.

In addition to substantially full-length polypeptides expressed by variant genes, the present invention includes

biologically active fragments of the polypeptides, or analogs thereof, including organic molecules which simulate the interactions of the peptides. Biologically active fragments include any portion of the full-length polypeptide which  
5 confers a biological function on the variant gene product, including ligand binding, and antibody binding. Ligand binding includes binding by nucleic acids, proteins or polypeptides, small biologically active molecules, or large cellular structures.

10 Polyclonal and/or monoclonal antibodies that specifically bind to variant gene products but not to corresponding prototypical gene products are also provided. Antibodies can be made by injecting mice or other animals with the variant gene product or synthetic peptide fragments  
15 thereof. Monoclonal antibodies are screened as are described, for example, in Harlow & Lane, *Antibodies, A Laboratory Manual*, Cold Spring Harbor Press, New York (1988); Goding, *Monoclonal antibodies, Principles and Practice* (2d ed.) Academic Press, New York (1986). Monoclonal antibodies are  
20 tested for specific immunoreactivity with a variant gene product and lack of immunoreactivity to the corresponding prototypical gene product. These antibodies are useful in diagnostic assays for detection of the variant form, or as an active ingredient in a pharmaceutical composition.

#### 25 V. Kits

The invention further provides kits comprising at least one allele-specific oligonucleotide as described above. Often, the kits contain one or more pairs of allele-specific oligonucleotides hybridizing to different forms of a  
30 polymorphism. In some kits, the allele-specific oligonucleotides are provided immobilized to a substrate. For example, the same substrate can comprise allele-specific oligonucleotide probes for detecting at least 10, 100 or all of the polymorphisms shown in Table 1. Optional additional  
35 components of the kit include, for example, restriction enzymes, reverse-transcriptase or polymerase, the substrate nucleoside triphosphates, means used to label (for example, an

avidin-enzyme conjugate and enzyme substrate and chromogen if the label is biotin), and the appropriate buffers for reverse transcription, PCR, or hybridization reactions. Usually, the kit also contains instructions for carrying out the methods.

5 VI. Computer Systems For Storing Polymorphism Data

Fig. 1A depicts a block diagram of a computer system  
10 suitable for implementing the present invention. Computer system 10 includes a bus 12 which interconnects major subsystems such as a central processor 14, a system memory 16  
10 (typically RAM), an input/output (I/O) controller 18, an external device such as a display screen 24 via a display adapter 26, serial ports 28 and 30, a keyboard 32, a fixed disk drive 34 via a storage interface 35 and a floppy disk drive 36 operative to receive a floppy disk 38, and a CD-ROM  
15 (or DVD-ROM) device 40 operative to receive a CD-ROM 42. Many other devices can be connected such as a user pointing device, e.g., a mouse 44 connected via serial port 28 and a network interface 46 connected via serial port 30.

Many other devices or subsystems (not shown) may be  
20 connected in a similar manner. Also, it is not necessary for all of the devices shown in Fig. 1A to be present to practice the present invention, as discussed below. The devices and subsystems may be interconnected in different ways from that shown in Fig. 1A. The operation of a computer system such as  
25 that shown in Fig. 1A is well known. Databases storing polymorphism information according to the present invention can be stored, e.g., in system memory 16 or on storage media such as fixed disk 34, floppy disk 38, or CD-ROM 42. An application program to access such databases can be operably  
30 disposed in system memory 16 or sorted on storage media such as fixed disk 34, floppy disk 38, or CD-ROM 42.

Fig. 1B depicts the interconnection of computer system  
10 to remote computers 48, 50, and 52. Fig. 1B depicts a network 54 interconnecting remote servers 48, 50, and 52.  
35 Network interface 46 provides the connection from client computer system 10 to network 54. Network 54 can be, e.g., the Internet. Protocols for exchanging data via the Internet

and other networks are well known. Information identifying the polymorphisms described herein can be transmitted across network 54 embedded in signals capable of traversing the physical media employed by network 54.

5 Information identifying polymorphisms shown in Table 1 is represented in records, which optionally, are subdivided into fields. Each record stores information relating to a different polymorphisms in Table 1. Collectively, the records can store information relating to all of the polymorphisms in  
10 Table 1, or any subset thereof, such as 5, 10, 50, or 100 polymorphisms from Table 1. In some databases, the information identifies a base occupying a polymorphic position and the location of the polymorphic position. The base can be represented as a single letter code (i.e., A, C, G or T/U)  
15 present in a polymorphic form other than that in the reference allele. Alternatively, the base occupying a polymorphic site can be represented in IUPAC ambiguity code as shown in Table 1. The location of a polymorphic site can be identified as its position within one of the sequences shown in Table 1.  
20 For example, in the first sequence shown in Table 1, the polymorphic site occupies the 16th base. The position can also be identified by reference to, for example, a chromosome, and distance from known markers within the chromosome. In other databases, information identifying a polymorphism  
25 contains sequences of 10-100 bases shown in Table 1 or the complements thereof, including a polymorphic site. Preferably, such information records at least 10, 15, 20, or 30 contiguous bases of sequences including a polymorphic site.

#### EXAMPLES

30 The polymorphisms shown in Table 1 were identified by resequencing of target sequences from eight unrelated individuals of diverse ethnic and geographic backgrounds by hybridization to probes immobilized to microfabricated arrays. The strategy and principles for design and use of such arrays  
35 are generally described in WO 95/11995. The strategy provides arrays of probes for analysis of target sequences showing a

high degree of sequence identity to the reference sequences of the fragments shown in Table 1, column 1. The reference sequences were sequence-tagged sites (STSs) developed in the course of the Human Genome Project (see, e.g., *Science* 270, 1945-1954 (1995); *Nature* 380, 152-154 (1996)). Most STS's  
5      ranged from 100 bp to 300 bp in size.

A typical probe array used in this analysis has two groups of four sets of probes that respectively tile both strands of a reference sequence. A first probe set comprises  
10     a plurality of probes exhibiting perfect complementarity with one of the reference sequences. Each probe in the first probe set has an interrogation position that corresponds to a nucleotide in the reference sequence. That is, the  
15     interrogation position is aligned with the corresponding nucleotide in the reference sequence, when the probe and reference sequence are aligned to maximize complementarity between the two. For each probe in the first set, there are  
20     three corresponding probes from three additional probe sets. Thus, there are four probes corresponding to each nucleotide in the reference sequence. The probes from the three  
25     additional probe sets are identical to the corresponding probe from the first probe set except at the interrogation position, which occurs in the same position in each of the four  
30     corresponding probes from the four probe sets, and is occupied by a different nucleotide in the four probe sets. In the  
35     present analysis, probes were 25 nucleotides long. Arrays tiled for multiple different reference sequences were included on the same substrate.

Multiple target sequences from an individual were  
30     amplified from human genomic DNA using primers for the fragments indicated in the listed Web sites. The amplified target sequences were fluorescently labelled during or after PCR. The labelled target sequences were hybridized with a  
35     substrate bearing immobilized arrays of probes. The amount of label bound to probes was measured. Analysis of the pattern of label revealed the nature and position of differences between the target and reference sequence. For example, comparison of the intensities of four corresponding probes

reveals the identity of a corresponding nucleotide in the target sequences aligned with the interrogation position of the probes. The corresponding nucleotide is the complement of the nucleotide occupying the interrogation position of the probe showing the highest intensity (see WO 95/11995). The existence of a polymorphism is also manifested by differences in normalized hybridization intensities of probes flanking the polymorphism when the probes hybridized to corresponding targets from different individuals. For example, relative loss of hybridization intensity in a "footprint" of probes flanking a polymorphism signals a difference between the target and reference (i.e., a polymorphism) (see EP 717,113, incorporated by reference in its entirety for all purposes). Additionally, hybridization intensities for corresponding targets from different individuals can be classified into groups or clusters suggested by the data, not defined *a priori*, such that isolates in a give cluster tend to be similar and isolates in different clusters tend to be dissimilar. See WO 97/29212 (incorporated by reference in its entirety for all purposes). Hybridizations to samples from different individuals were performed separately. Table 1 summarizes the data obtained for target sequences in comparison with a reference sequence for the eight individuals tested.

From the foregoing, it is apparent that the invention includes a number of general uses that can be expressed concisely as follows. The invention provides for the use of any of the nucleic acid segments described above in the diagnosis or monitoring of diseases, such as cancer, inflammation, heart disease, diseases of the CNS, and susceptibility to infection by microorganisms. The invention further provides for the use of any of the nucleic acid segments in the manufacture of a medicament for the treatment or prophylaxis of such diseases. The invention further provides for the use of any of the DNA segments as a pharmaceutical.

All publications and patent applications cited above are incorporated by reference in their entirety for all

purposes to the same extent as if each individual publication or patent application were specifically and individually indicated to be so incorporated by reference. Although the present invention has been described in some detail by way of  
5 — illustration and example for purposes of clarity and understanding, it will be apparent that certain changes and modifications may be practiced within the scope of the appended claims.



WHAT IS CLAIMED IS:

- 1           1     A nucleic acid segment of between 10 and 100  
2 bases from a fragment shown in Table 1 including a polymorphic  
3 site, or the complement of the segment.
- 1           2.     The nucleic acid segment of claim 1 that is  
2 DNA.
- 1           3.     The nucleic acid segment of claim 1 that is RNA.
- 1           4.     The segment of claim 1 that is less than 50  
2 bases.
- 1           5.     The segment of claim 1 that is less than 20  
2 bases.
- 1           6.     The segment of claim 1, wherein the fragment is  
2 WI-14263 and the polymorphic site is at position 49.
- 1           7.     The segment of claim 1, wherein the polymorphic  
2 site is diallelic.
- 1           8.     The segment of claim 1, wherein the polymorphic  
2 form occupying the polymorphic site is the reference base for  
3 the fragment listed in Table 1, column 3.
- 1           9.     The segment of claim 1, wherein the polymorphic  
2 form occupying the polymorphic site is an alternative form for  
3 the fragment listed in Table 1, column 5.
- 1           10.    An allele-specific oligonucleotide that  
2 hybridizes to a segment of a fragment shown in Table 1, column  
3 8 or its complement.
- 1           11.    The allele-specific oligonucleotide of claim 10  
2 that is probe.

1           12. The allele-specific oligonucleotide of claim 10,  
2 wherein a central position of the probe aligns with the  
3 polymorphic site of the fragment.

1           13. The allele-specific oligonucleotide of claim 10  
2 that is a primer.

1           14. The allele-specific oligonucleotide of claim 13,  
2 wherein the 3' end of the primer aligns with the polymorphic  
3 site of the fragment.

1           15. An isolated nucleic acid comprising a sequence of  
2 Table 1, column 8 or the complement thereof, wherein the  
3 polymorphic site within the sequence or complement is occupied  
4 by a base other than the reference base show in Table 1,  
5 column 3.

1           16. A method of analyzing a nucleic acid, comprising:  
2 obtaining the nucleic acid from an individual; and  
3 determining a base occupying any one of the polymorphic  
4 sites shown in Table 1.

1           17. The method of claim 16, wherein the determining  
2 comprises determining a set of bases occupying a set of the  
3 polymorphic sites shown in Table 1.

1           18. The method of claim 16, wherein the nucleic acid  
2 is obtained from a plurality of individuals, and a base  
3 occupying one of the polymorphic positions is determined in  
4 each of the individuals, and the method further comprising  
5 testing each individual for the presence of a disease  
6 phenotype, and correlating the presence of the disease  
7 phenotype with the base.

8           19. A computer-readable storage medium for storing  
9 data for access by an application program being executed on a  
10 data processing system, comprising:

11                   a data structure stored in the computer-readable  
12 storage medium, the data structure including information  
13 resident in a database used by the application program and  
14 including:

15                   a plurality of records, each record of the  
16 plurality comprising information identifying a polymorphisms  
17 shown in Table 1.

18           20. The computer-readable storage medium of claim 19,  
19 wherein each record has a field identifying a base occupying a  
20 polymorphic site and a location of the polymorphic site.

21           21. The computer-readable storage medium of claim 19,  
22 wherein each record identifies a nucleic acid segment of  
23 between 10 and 100 bases from a fragment shown in Table 1  
24 including a polymorphic site, or the complement of the  
25 segment.

26           22. The computer-readable storage medium of claim 19,  
27 comprising at least 10 records, each record comprising  
28 information identifying a different polymorphism shown in  
29 Table 1.

30           23. The computer-readable storage medium of claim 19,  
31 comprising at least 100 records, each record comprising  
32 information identifying a different polymorphisms shown in  
33 Table 1.

34           24. A signal carrying data for access by an  
35 application program being executed on a data processing  
36 system, comprising:  
37                   a data structure encoded in the signal, said data  
38 structure including information resident in a database used by  
39 the application program and including:  
40                   a plurality of records, each record of the plurality  
41 comprising information identifying a polymorphism shown in  
42 Table 1.

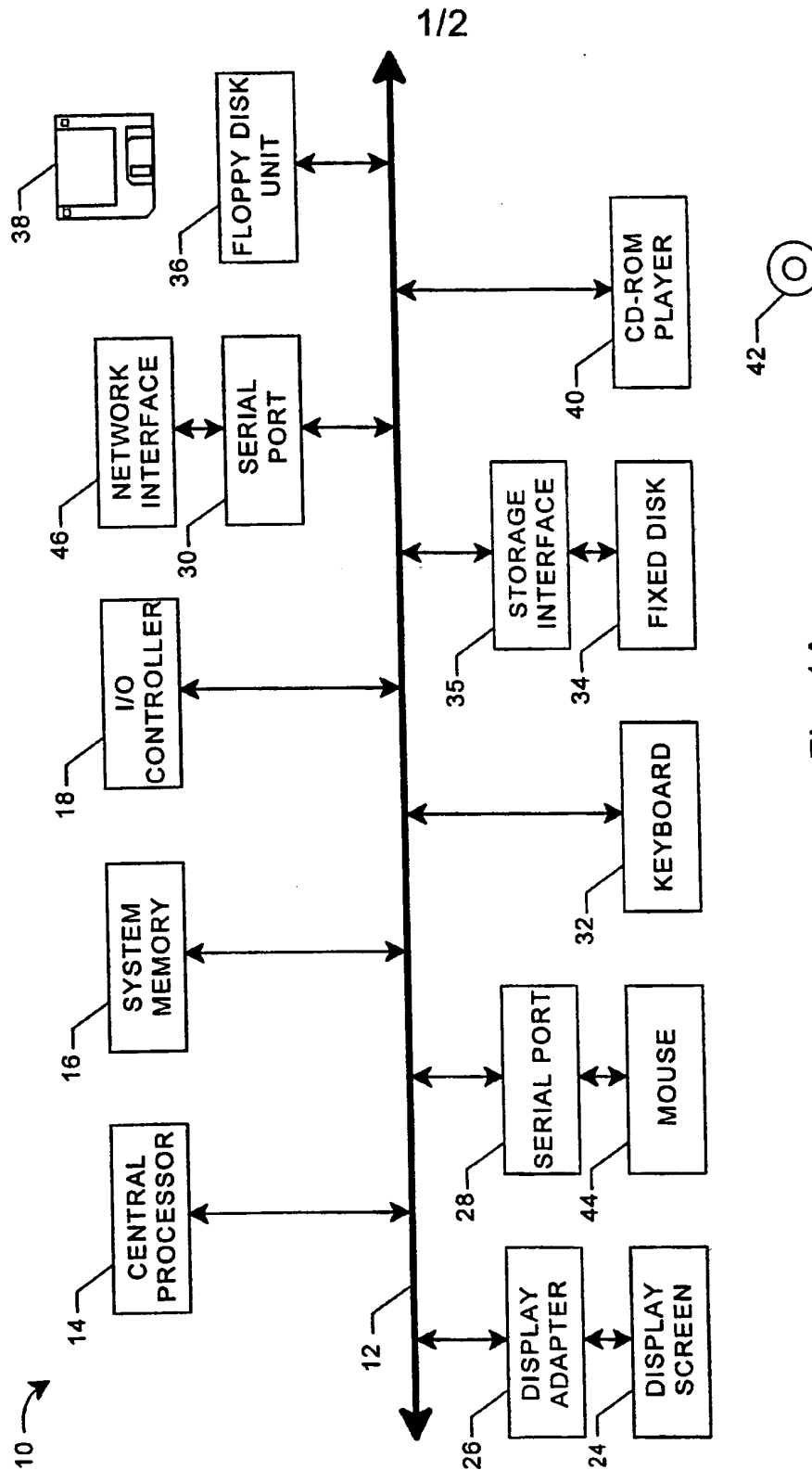


Fig. 1A

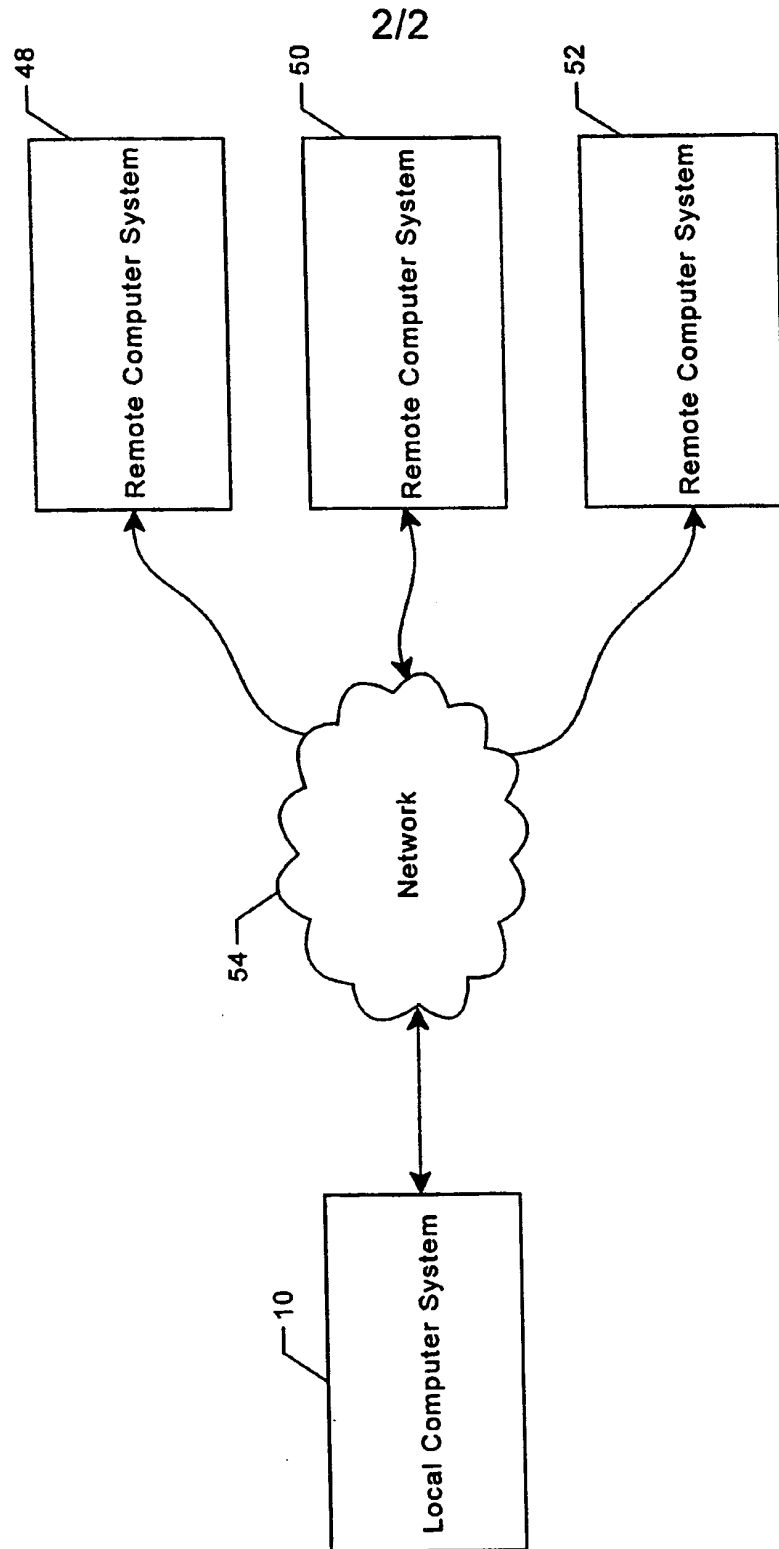


Fig. 1B